



Dataset Article

The Global Wheat Full Semantic Organ Segmentation (GWFSS) dataset[☆]

Zijian Wang^{l,1}, Radek Zenkl^{f,1}, Latifa Greche^{k,1}, Benoit De Solan^a,
 Lucas Bernigaud Samatan^{a,2}, Safaa Ouahid^{b,o}, Andrea Visioni^b, Carlos A. Robles-Zazueta^{c,3},
 Francisco Pinto^{c,4}, Ivan Perez-Olivera^{c,5}, Matthew P. Reynolds^c, Chen Zhu^d, Shouyang Liu^d,
 Marie-Pia D'argaignon^{e,6}, Raul Lopez-Lozano^e, Marie Weiss^e, Afef Marzougui^f, Lukas Roth^f,
 Sébastien Dandrifosse^{g,7}, Alexis Carlier^{g,8}, Benjamin Dumont^g, Benoît Mercatoris^g,
 Javier Fernandez^h, Scott Chapman^h, Keyhan Najafianⁱ, Ian Stavnessⁱ, Haozhou Wang^j,
 Wei Guo^j, Nicolas Virlet^k, Malcolm J. Hawkesford^k, Zhi Chen^l, Etienne David^m, Joss Gilletⁿ,
 Kamran Irfanⁿ, Alexis Comarⁿ, Andreas Hund^{f,*}

^a Arvalis, France

^b International Center for Agricultural Research in the Dry Areas, Rabat, Morocco

^c Global Wheat Program, International Maize and Wheat Improvement Center, C.P. 56237, El Batan, Texcoco, Mexico

^d Nanjing Agricultural University, China

^e EMMAH, UMR1114, LPA CAPTE, INRAE, Avignon, France

^f Department of Environmental Systems Science, ETH Zurich, 8092, Zurich, Switzerland

^g Plant Sciences & Biosystems Dynamics and Exchanges, TERRA Teaching and Research Centre, Gembloux Agro-Bio Tech, University of Liège, Gembloux, Belgium

^h School of Agriculture and Food Sustainability, The University of Queensland, Brisbane, Australia

ⁱ Department of Computer Science, University of Saskatchewan, Saskatoon, Canada

^j Graduate School of Agricultural and Life Sciences, The University of Tokyo, Tokyo, Japan

^k Sustainable Soils and Crops, Rothamsted Research, West Common, Harpenden, AL5 2JQ, UK

^l School of Electrical Engineering and Computer Science, The University of Queensland, Brisbane, Australia

^m neoBloom, Munstermannskamp 1, 21335 Lüneburg, Germany

ⁿ HIPHEN SAS, 120 rue Jean Dausset, 84140 Avignon, France

^o Programa de Doctorado de Ingeniería Agraria, Alimentaria, Forestal y del Desarrollo Rural Sostenible, Universidad de Córdoba, Córdoba, Spain

ARTICLE INFO

Keywords:

Wheat organ segmentation
 Field phenomics
 High-throughput phenotyping
 Breeding

ABSTRACT

Computer vision is increasingly used in farmers' fields and agricultural experiments to quantify important traits. Imaging setups with a sub-millimeter ground sampling distance enable the detection and tracking of plant features, including size, shape, and colour. Although today's AI-driven foundation models segment almost any object in an image, they still fail for complex plant canopies. To improve model performance, the global wheat dataset consortium assembled a diverse set of images from experiments around the globe. After the head detection dataset (GWHHD), the new dataset targets a full semantic segmentation (GWFSS) of organs (leaves, stems and spikes) covering all developmental stages. Images were collected by 11 institutions using a wide range of imaging setups. Two datasets are provided: i) a set of 1096 diverse images in which all organs were labelled at

[☆] Institutions contributing to the dataset are listed in alphabetic order with the same numbering as in Section 2.1 followed by institution without image contribution. The group leaders within institution are listed last. Exception: Lead authors are listed first; contributing author is listed last.

* Corresponding author.

E-mail address: andreas.hund@usys.ethz.ch (A. Hund).

¹ These authors contributed equally to this work.

² Current address: University of Toulouse, INRAE, UMR DYNAPOR, 31326 Castanet-Tolosan, France.

³ Current address: Department of Plant Breeding, Hochschule Geisenheim University, 65366 Geisenheim, Germany.

⁴ Current address: Centre for Crop Systems Analysis, Wageningen University & Research, Bornsesteeg 48, Building 109, 6708 PE Wageningen, The Netherlands.

⁵ Current address: Department of Agricultural and Biosystems Engineering, South Dakota State University, P. O. Box 57007, Brookings, SD, USA.

⁶ Current address: Limagrain Chappes - Centre de Recherche; 28 Rte d'Ennezat 63720 Clermont-Ferrant.

⁷ Current address: Walloon Agricultural Research Centre, Gembloux, Belgium.

⁸ Current address: Osiris Agriculture, France.

<https://doi.org/10.1016/j.plaphe.2025.100084>

Received 11 March 2025; Received in revised form 27 May 2025; Accepted 8 June 2025

Available online 6 August 2025

2643-6515/© 2025 The Authors. Published by Elsevier B.V. on behalf of Nanjing Agricultural University. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

the pixel level, and (ii) a dataset of 52,078 images without annotations available for additional training. The labelled set was used to train segmentation models based on DeepLabV3Plus and Segformer. Our Segformer model performed slightly better than DeepLabV3Plus with a mIOU for leaves and spikes of ca. 90 %. However, the precision for stems with 54 % was rather lower. The major advantages over published models are: i) the exclusion of weeds from the wheat canopy, ii) the detection of all wheat features including necrotic and senescent tissues and its separation from crop residues. This facilitates further development in classifying healthy vs. unhealthy tissue to address the increasing need for accurate quantification of senescence and diseases in wheat canopies.

1. Introduction

Wheat is one of the most important crops in the world, providing 18 % of calorie intake and 19 % of protein intake globally [1]. While the area of farmland used to grow wheat has remained stable, global wheat yields have quadrupled since 1960, mainly due to technical innovation, such as the widespread use of N fertilisers, and the breeding of modern wheat varieties [1]. However, the rate of yield gain has stagnated or even decreased in the last two decades in different regions of the world [2], imposing a challenge to fulfilling the projected demands for wheat production in the future [3]. In addition, climate change-induced stresses [1], plant disease adaptation, and pest outbreaks, combined with the need for more efficient crops that require less input in terms of fertilisers, water, and pesticides place additional constraints on wheat production. Increasing wheat yield is a multi-faceted problem which involves genetic, physiologic, and agronomic improvement to enhance resource-use efficiency. Novel phenotyping approaches provide advanced tools and methodologies to enhance wheat management, optimise breeding, and achieve efficient resource utilization. These methodologies are key to enabling precise repeatable measurements in agricultural fields and research networks across the globe, as highlighted in a survey on field-based phenotyping in Europe [4]. In their survey, Morisse et al. [4] emphasise the capability of these platforms to enable the rapid collection of datasets at breeding scale, including hundreds of plots with various genotypes and treatments. Moreover, data can be collected during the entire crop growth cycle, enabling tracking of wheat responses to biotic and abiotic stresses or management practices. Although a wide range of possible sensors are available, including 3D scanners [5] to characterise canopy architecture or hyperspectral imaging [6,7] to monitor plant health and productivity, we focus here on high spatial-resolution imaging with red, green and blue (RGB) spectra. Such imaging provides a broad range of phenotyping capabilities at relatively low equipment costs and high spatial resolution. In a comprehensive review on translating high-throughput phenotyping (HTP) into genetic gain, Araus et al. [8], posed a key question: “Will low-cost HTP tools be adopted regularly by breeders in the next decades? If so, are RGB cameras, mobile apps, and drones the natural candidates?” Classical RGB cameras can be mounted on handheld devices [9], ground-based vehicles [8,10], gantries [11], or drones [12], making them adaptable to various scales and platforms. The advantage of RGB sensors is their relatively high spatial resolution and low cost compared to other HTP tools [13]. This enables the capture of sufficient details to separate plant features from complex canopies. However, robust feature extraction requires algorithms capable of extracting information from images taken under a wide range of conditions. The development of such algorithms demands a sufficient number of well-annotated training images capturing this diversity.

Several pioneer datasets published related wheat-linked tasks in field conditions focusing either on (i) spike detection and quantification with the SPIKE [14] and GHWD dataset [15,16], and a dataset provided by Madec et al. [17], (ii) vegetation segmentation on wheat only with EWS dataset [18] or on multispecies including wheat with SegVeg and VegAnn [9,19] or (iii) disease detection on wheat with the NWRD dataset [20], the wheat leaf dataset for strip rust and septoria [21], the CDTs for strip rust [22], the wheat nitrogen deficiency and leaf rust

image dataset [23] and the CGIAR Computer Vision for Crop Disease for stem and leaf rust [24]. Despite these contributions, existing datasets are often limited in terms of geographic diversity, genotype variations, and growth stages, which may limit the generalisation power of models trained on these datasets. Furthermore, only a few datasets are collaborative, involving multiple countries and institutions [9,15,16,19]. To address this limitation, the collaborative Global Wheat Dataset Consortium was established. This consortium aims to aggregate datasets from multiple institutions, make them publicly available, and provide the necessary data to develop robust algorithms. In previous GWHD editions 2020–2021 [15,16], we released a large dataset with a total of 6515 high-resolution RGB images, containing annotations for 275,187 labelled wheat heads. The early design of GWHD focused on providing bounding-box annotations for wheat heads, including images of different genotypes captured under varying environmental, management and growth conditions. GWHD played a pivotal role in various research studies, particularly in developing and benchmarking wheat head detection and counting methods using supervised models [25–31], semi-self-supervised models [32,33] and self-supervised models [34], generating a reference dataset to improve deep learning model performance [35–37] and improving head count models in dense plots [38]. These robust detection models allow the counting of wheat heads per unit area [29,39], as long as the footprint of the image at the top of the canopy is known. In addition, the dynamics of head counts may be used to approximate heading dates [40]. The flowering date of wheat is often approximated by the heading date because it is easier to assess head emergence or presence than anther extrusion which is more affected by time of day, wind conditions and operator experience [40,41]. Thus, heading is the most widely assessed trait related to cereal phenology and an important trait to understand the effect of environmental stresses, such as heat and drought, on grain yield [42].

While counting wheat heads is important as it is one of the yield components, it is not the only targeted trait. Wheat continuously adjusts its yield potential during the entire vegetation period. Low germination or plant damage due to winter kill is compensated for by increased tillering, while during stem elongation, excessive tillering is compensated by tiller abortion. Later in the season, the different organs can undergo different senescence dynamics, as demonstrated by Anderegg et al. [43]. A deeper understanding of how yield is formed throughout the growth season will benefit from a non-destructive assessment of its components. The imaging and semantic segmentation of all the plant organs visible in the image, from emergence to maturity, have great potential to shed light on the yield formation process. Examples of targeted traits are seedling count [44,45], canopy cover [18,19,46,47], biomass estimation [48] and leaf area index (LAI) [49,50]. In most cases, segmentation of wheat canopies from other background features, such as soil or weeds, is required. The comparably simple task of segmenting canopies from the soil background was previously solved by manual adjustments [51], using automatic threshold methods, such as the Otsu algorithm [52,53] or shallow machine learning [13,54,55].

The above-mentioned approaches based on colour information at the pixel-level, have the shortcoming that they cannot take context into account. Modern deep learning methodologies enable the learning of contextual information. This requires human-annotated training data to supervise feature detection in complex images containing many plants

growing together in a canopy. The different plant organs in such canopies are not simply green but might have different shades of green or yellow due to senescence, chlorosis, or necrotic tissues. Necrotic plant tissues may have a similar brown colour as crop residues and can only be segmented in RGB images based on context. Green canopy segmentation using deep learning models has become a standard procedure, with training data sourced from a wide range of crops, such as the VegAnn dataset [9] combining 3775 RGB images of 12 different crops. With the existence of large datasets and the advancements in computer vision, new possibilities of data processing and feature extraction have been unlocked through the utilization of data-driven deep learning approaches. Existing deep learning-based segmentation is mainly based on encoder-decoder architecture like DeeplabV3Plus [56–60] and Atrous convolution [61,62] architecture, with application-specific adaptation. More recently, transformer-based segmentation models (e.g., SegFormer [63]) have gained attention due to their ability to capture long-range dependencies and global context effectively. These models show promise in addressing complex segmentation tasks, offering improved performance and adaptability in agricultural applications. Despite their potential to achieve promising performance, transformer-based models require relatively larger datasets that are currently lacking in agricultural domains [64].

Within the wheat canopy, segmentation of organs including leaves, stems, and heads is required, followed by the extraction of relevant phenotypes from each organ. Many wheat researchers have collected their own wheat datasets or made their own annotations from existing ones to achieve semantic segmentation of heads [35,59,65–67], spikelets [68], grains [69], stems and foliage [43] infected [20,70–72] and senescent [19] tissues, or a combination of disease and senescence [73]. The data provided in these publications largely advance wheat phenotyping at the organ level, offering tools for detailed studies of yield components such as spike number, spikelets per spike, spikelet size, leaf disease resistance, and senescence dynamics. However, there is a lack of integrated datasets enabling simultaneous segmentation of all wheat organs (leaves, spikes, stems) from soil, crop residues, weeds and other background elements. Moreover, wheat has a complex canopy due to its high planting density, strong development of tillers (lateral shoots), thin stems, overlapping leaves and occluded organs. Variations in appearance caused by growth stages, lighting conditions, wind patterns and imaging angles make it a challenging plant species to phenotype. To advance our capabilities beyond the wheat head-centred GWHD dataset 2021–2022, we assembled the Global Wheat Full Semantic Segmentation (GWFS) dataset collected by different phenotyping platforms from 11 institutes and universities across the globe under various light and weather conditions. This diversity ensures that the dataset addresses the extensive requirements of wheat phenotyping across a range of genetic backgrounds, environments, and management practices throughout the growing season. The images were collected with an average ground sampling distance (GSD) between 0.09 and 0.71 mm per pixel. This is significant for accurately capturing organ features to allow precise differentiation and measurement of these smaller structures, rather than focusing solely on canopy-level traits. Our contribution can be summarised as follows:

1. A *full* GWFS dataset comprising 52,078 RGB images without labels is available in ETH research collection.
2. An *annotated* GWFS dataset providing 1096 pixel-level annotations (masks) for the following classes: leaves, stems, heads, and background.
3. The results of two state-of-the-art segmentation models, DeepLabV3Plus and Segformer, fine-tuned on the full dataset. The models were trained as a baseline performance benchmark for organ segmentation.

2. Material and methods

2.1. Field experiment

The dataset includes images from field experiments conducted by 11 institutions worldwide, as detailed in Table 1. Wheat plots at 67 different field sites were imaged using proximal RGB imaging setups throughout the growing seasons (Fig. 1). The experiments cover a wide range of planting densities, agronomic inputs, environmental conditions, as well as disease and weed pressures. Thus, the GWFS dataset spans diverse agroclimatic zones and management practices. The imaging setups used by the 11 institutions and the data sets derived are described in detail in Table S1 and Table 2, respectively. Additional information related to the datasets is given as follows with institutions listed in alphabetic order:

1. Arvalis

ARVALIS_1–200: The 200 subsets were acquired in 2022 and 2023 in a network of 18 sites representing the main agroclimatic zones and the most common practices in France. The trials cover different themes: evaluation of different wheat genotypes at diverse nitrogen fertilisation regimes, management methods, diseases, pests and water stress. Some trials comply with the specifications of organic farming or include wheat in combination with other species. Seed densities are a typical common practice in France. Images at most sites have been collected with the LITERAL [74], a handheld system with high-resolution cameras working without flash; The images at the location GREOUX were acquired by PHENOMOBILE [9], an autonomous robot equipped with industrial cameras and flashlights.

2. International Center for Agricultural Research in the Dry Areas (ICARDA)

ICARDA_MCH_2023: Data collected at Merchouche Station (the main ICARDA experimental station near Rabat, Morocco), under drought conditions. This set includes 960 entries from stage 2 of the durum wheat breeding program. Three acquisitions were taken with ICARDA's PHENOBUGGY (equipped with RGB camera, multispectral and LiDAR) in March 2023. *ICARDA_MCH_2024*: Data collected at Merchouche, also during a quite dry year. This set comprise 3 different trials (5 acquisitions each): 1) CWR panel as part of the Crop Trust project BOLD representing 60 elite durum wheat lines obtained from crosses with crop wild relatives; 2) CEREALMED including 288 entries of the Durum

Table 1

Summary of Institutions, the number of unique geolocations, the number of unique Image Setups, and the number of images per subset in the full dataset. A full description of the datasets (GWFS_v1.0_subsets.csv) is available in the ETH research collection referenced in the Data Availability section.

	Institution	Country	Geolocations (by lat/long)	Image Setups	Images in full set
1	ARVALIS	France	17	3	4134
2	ICARDA	Morocco	1	1	5088
3	CIMMYT-CENEB	Mexico	1	1	5773
4	NJAU	China	4	2	5008
5	INRAE	France	3	2	4102
6	ETHZ	Switzerland	2	3	7175
7	ULIEGE-CRAW	Belgium	7	3	5199
8	UQ	Australia	28	2	4265
9	UTokyo	Japan	1	1	6130
10	USask	Canada	1	1	200
11	RRes	UK	1	3	5004

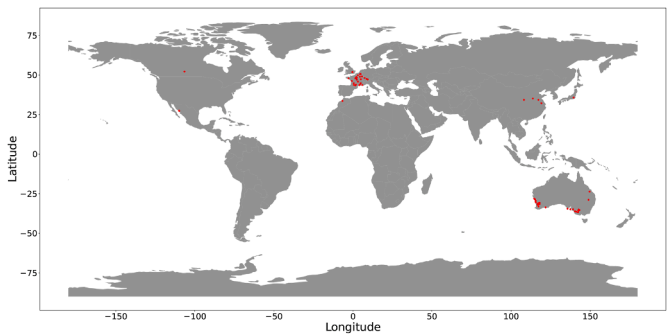


Fig. 1. An overview of the location of all trials included in GWFS.

- Global Panel of landraces, modern and old varieties mainly from Mediterranean countries; 3) a root rot trial of 24 elite lines.
3. **International Maize and Wheat Improvement Center (CIMMYT)**
- CIMMYT-CENEB_1–7*: The CIMMYT dataset includes images collected for 319 spring bread wheat genotypes, consisting of elite, pre-breeding and exotic germplasm phenotyped during 2020 and 2021 field seasons in Campo Experimental Norman E. Borlaug (CENEB) in Ciudad Obregon, Sonora, Mexico. The genotypes were imaged from heading to maturity under irrigated, drought and terminal heat stress conditions in the field.
4. **Nanjing Agricultural University (NJAU)**
- The NJAU datasets were collected from field trials in different regions of China. Phenotyping data for *NJAU_1–NJAU_4* were collected using PhenoArm, a portable handheld imaging platform with two high-resolution cameras. *NJAU_5* used Phenotypette, a pushcart platform integrating LiDAR, multispectral, and RGB cameras. The pushcart was manually operated at a controlled speed and equipped with RTK-GPS for automated data collection. *NJAU_1–2*: Experiments in Jurong and Xuzhou during 2020–2021, with 5 wheat cultivars under 3 nitrogen levels.

- NJAU_3*: Trial in Xinxiang from 2022 to 2023 including 565 wheat cultivars, covering both introduced and domestic cultivars since 1950. Five cultivars were replicated 16 times, while the remaining 560 cultivars had no replication. A total of 640 plots were established, with fertilisation and irrigation managed according to local practices.
- NJAU_4–5*: Trials in Yangling from 2021–2024 with the same cultivars as *NJAU_3*.
5. **National Research Institute for Agriculture, Food and Environment (INRAE)**
- The INRAE dataset was acquired in the frame of the FFAST project (French National Grant ANR-21-CE45-0037). The dataset includes images taken in field trials at three INRAE experimental sites UE APC at Auzerville (AUZ), UE DiaScope at Mauguio (MAU) and UE PHACC at Clermont-Ferrand (CLE) in the years 2021, 2022 and 2023. All pictures were taken using the Phenomobile V2 ground robot (<https://hal.inrae.fr/hal-03646863>), equipped with RGB cameras looking at nadir and at 45°. Images were taken in active illumination conditions (flashes). The trials consisted of 10 French elite cultivars grown under 4 treatments (depending on the site: irrigation, sowing date and seed density).
6. **Swiss Federal Institute of Technology Zurich; ETH Zurich (ETHZ)**
- ETHZ_01*: Images from the ‘field phenotyping platform’ (FIP) at ETH Zurich in Eschikon [75]. The site covers typical climatic conditions of the Swiss Plateau. About 350 wheat varieties are monitored at least once per week to relate growth patterns to causal environmental factors. The set is available at [76]. *ETHZ_02*: Organic farming conditions at 981 m altitude, long snow cover and 2084 mm annual precipitation peaking in summer. Images show damage in spring caused by snow mould (Microdochium nivale). The high precipitation fostered lodging and diseases in the summer. There was high weed pressure.
7. **University of Liège and Walloon Agricultural Research Center (ULIEGE-CRAW)**
- ULIEGE-CRA-W_01–18*: Images were acquired in winter wheat

Table 2
Imaging Setup Details. A full description of the imaging setups (GWFS_v1.0_imaging_setups.csv) is available in the ETH research collection referenced in the Data Availability section.

Imaging Setup	Vector	Camera Model	Viewing Angle	Focal Length (mm)	Sensor Resolution (pixel)	Field of View Horizontal (°)	Field of View Vertical (°)	Distance to Ground (m)	GSD (mm/px)
LITERAL1.0_0	Handheld	Sony RX0	0	7.7	4800 × 3200	73	52	1.8	0.55
LITERAL1.0_45	Handheld	Sony RX0	45	7.7	4801 × 3200	73	52	1.2	0.32
FIP1.0	Gantry	Canon EOS 5D Mark II	0	35	5616 × 3744	54.43	37.85	3	0.5495
GO1.0M_0	Cart	JAI GO-5000C-USB	0	16	2560 × 2048	44.3	33.6	1 (to canopy)	0.3125
GO1.0M_30	Cart	JAI GO-5000C-USB	30	16	2560 × 2048	44.3	33.6	1 (to canopy)	0.3608
GO1.6M_0	Cart	JAI GO-5000C-USB	0	16	2560 × 2048	44.3	33.6	1.6 (to canopy)	0.5
Rres_GT3300_top	Gantry	Prosilica GT3300C	0	50	3296 × 2472	38	26	2.8	0.12
Rres_GT3300_south	Gantry	Prosilica GT3300C	30	50	3296 × 2472	38	26	1.75	0.09
Rres_GT3300_north	Gantry	Prosilica GT3300C	30	50	3296 × 2472	38	26	1.75	0.09
PhenoArm	Handheld	Sony RX0	0, 45	7.7	4800 × 3200	81.2	59.5	2	0.7143
Phenotypette	Cart	Sony RX0	0, 45	7.7	4800 × 3200	81.2	59.5	2	0.7143
Low-cost phenomobile	Cart	Canon EOS 600D	0	55	1920 × 1080	23	15	2.4	0.046
PHENOMOBILE	Ground Vehicle	Baumer HXG40	0	25	2040 × 2040	28	28	1.8	0.43
Phenobuggy	Tractor Fobro	Baumer VCXG-124C	0	25	4096 × 3000	32	24	1.8–2.4	0.24–0.33
UFPS	Ground Vehicle	FLIR Chameleon3 USB3	0	16	2448 × 2048	26.4	19.8	2	0.45

trials in the Hesbaye area (Belgium) between 2018 and 2022. The 18 subsets detail the differences between the trials. Images cover mainly nitrogen fertilisation trials and nitrogen fertilisation \times fungicide trials [10]. Images also cover N, P and K fertilisation trials (ULIEGE-CRA-W_04) and drought experiments (ULIEGE-CRA-W_12, 18). The set also contains sample images from dense time series of the same plots recorded in 15-min intervals (ULIEGE-CRA-W_11). One series contains green reference spheres used as control points in thermal images acquired in addition to RGB images (ULIEGE-CRA-W_17).

8. The University of Queensland (UQ)

UQ_1–29: The 29 subsets detail the differences between trials in which the images were acquired. Images are collected in the 2020 and 2021 National Variety Trials. Differences include variations in geolocation, genotype, and growth stage across Australia. The photographs were taken in 2020 and 2021, using smartphone cameras from a top-down perspective at about 0.5–1.5 m above the canopy.

9. University of Saskatchewan (USask)

The USask dataset was collected from wheat phenotyping field trials at the Kernen Crop Research Farm in Saskatchewan, Canada in 2019. The images comprise a single field trial with 32 diverse wheat cultivars at the heading stage. Images were collected with the University of Saskatchewan Field Phenotyping System (UFPS), a custom-built, self-propelled ground vehicle equipped with a range of imaging instrumentation, RTK-GPS, and on-board data processing.

10. University of Tokyo (UTokyo)

The UTokyo dataset was collected from wheat phenotyping field trials at the Institute for Sustainable Agro-ecosystem Services (ISAS) in Tokyo, Japan, in the 2014–2015 season. A Field Server system [77] collected images of five genotypes through the whole growth stage. The camera module of the system is based on a digital single-lens reflex (DSLR) camera, the Canon EOS Kiss X5 camera, with an EF-S18-55 mm lens (Canon Inc., Tokyo) that provides high-quality and high-resolution (18 megapixels) image data. A preprogrammed microcontroller board controls the power and shutter of the camera automatically.

11. Rothamsted Research (RRes)

The Rothamsted dataset includes images collected for 391 wheat genotypes, captured throughout the growth cycle, from tillering to maturity, using the LemnaTec Field Scanalyzer [11]. The NIT subsets relate to the evaluation of four commercial variety growing supply with six levels of nitrogen input over two years (2019 and 2021). Images of the NIT subsets were captured from three angles: 30° north, 0° top, and 30° south, providing comprehensive spatial coverage. The PxCS and the PxG subsets provide images from two mapping populations that were planted in 2019 and 2021, respectively. The populations displayed a large range of variation in terms of phenology and height.

2.2. Image acquisition

The imaging setups consist of a vector and a camera ranged from hand-held over manual push-cart to fully automated rovers, gantry systems or a cable-suspended system mounted on poles (Table 2). Images were acquired with RGB cameras of at least 1920 \times 1080 pixel sensor resolution, oriented from nadir (0°) to a 45° viewing angle. All carriers positioned the camera between 3 m and 1 m above the ground, leading to ground sampling distances between 0.09 and 0.71 mm.

2.3. Data selection for the annotation pool

For data selection, we entrusted expert judgment. To assemble a diverse set of images for annotation, each participating institution was asked to provide approximately 5000 images encompassing different

phenological stages for the dataset and a diverse subset of 200 for the annotation pool. The selection of images prioritised diversity across key factors, including variations in phenological stages, geographic locations, cultivars, and imaging conditions. Additionally, institutions were encouraged to include treatments, such as varying nitrogen levels, irrigation, or other input variations, aiming to encompass a broad spectrum of scenarios in the dataset. From the annotation pool, a total of 1096 images were selected through a stratified approach, ensuring proportional representation across contributing institutions. Specifically, the selected image set comprises 110 images from each of the seven institutions (*i.e.*, INRAE, ETHZ, USASK, Arvalis, RRES, NJAU, and UQ), 109 images from two institutions (*i.e.*, ULiege and Utokeyo), and 108 images from CIMMYT.

The selection process was driven by a joint consideration of feature geometry and institutional balance. Specifically, image features were extracted using a ResNet model [78] pre-trained on ImageNet. A k-means clustering algorithm was then applied to group the images based on their feature similarity. To ensure representative sampling, images closest to the cluster centres were selected, while also maintaining a balanced distribution across institutions. This process was designed to maximize the uniqueness of the selected images based on their embedding distributions. As a final data preprocessing step, all images were standardised to ensure consistent resolution and comparable ground sampling distances. This was accomplished by applying a centre crop to achieve a resolution of 512 \times 512 pixels for most datasets. An exception was made for data contributed by UTokyo, where a resolution of 1024 \times 1024 pixels was used to accommodate the more detailed ground sampling distance. Fig. 2 illustrates the images sampled from our proposed dataset, showcasing its diversity.

2.4. Labeling

2.4.1. Targeted wheat entities

The annotation process was carried out centrally by expert annotators using the Darwin annotation tool provided by V7 Darwin.⁹ During annotation, temporary adjustments of brightness and contrast were done to enhance the distinction among features. The annotation process and quality control were handled by HIPHEN. In case of annotation mistakes, images were sent back to the annotators with respective instructions. In these cases, segmentation masks were modified using the brush and eraser tools. The global wheat experts team for labelling reviewed and resolved the unclear cases as needed. We refer readers to the appendix for the detailed GWFFS labelling guide. Initially, a small set of image tags was assigned to enhance understanding of image content and quality (see Table 3). Subsequently, pixel-level annotations were performed for the following classes: head, leaf, stem and background. As a reference for tissue types, we largely used the BRENDA Tissue Ontology (BTO¹⁰) retrieved in the EMBL-EBI Ontology Lookup Service.¹¹ The targeted entities were i) “heads” defined as spike (BTO_0001278) excluding awns (BTO_0005641), ii) “leaves” defined as the leaf lamina (BTO_0000719) including ligule, and iii) “stems” (BTO_0001300) including the surrounding leaf sheath (BTO_0005094). The peduncle was not labelled separately but included in “stem” (Fig. 3). The peduncle (PO_0009053) is the shoot axis that extends from the last foliage leaf on a stem (*i.e.* the flag leaf) until the next distal node (*i.e.* the basal end of spike).

Consequently, the wheat stem labelling included leaf sheaths, peduncles and bare stems (*e.g.* towards the end of the growing season). This decision considers that it is difficult to separate the different classes in complex images. The clear identification of the peduncle requires a visible ear (indicated by “P” in Fig. 3 d); the clear identification of a leaf

⁹ <https://darwin.v7labs.com>.

¹⁰ <https://bioportal.bioontology.org/ontologies/BTO>.

¹¹ <https://www.ebi.ac.uk/ols4>.

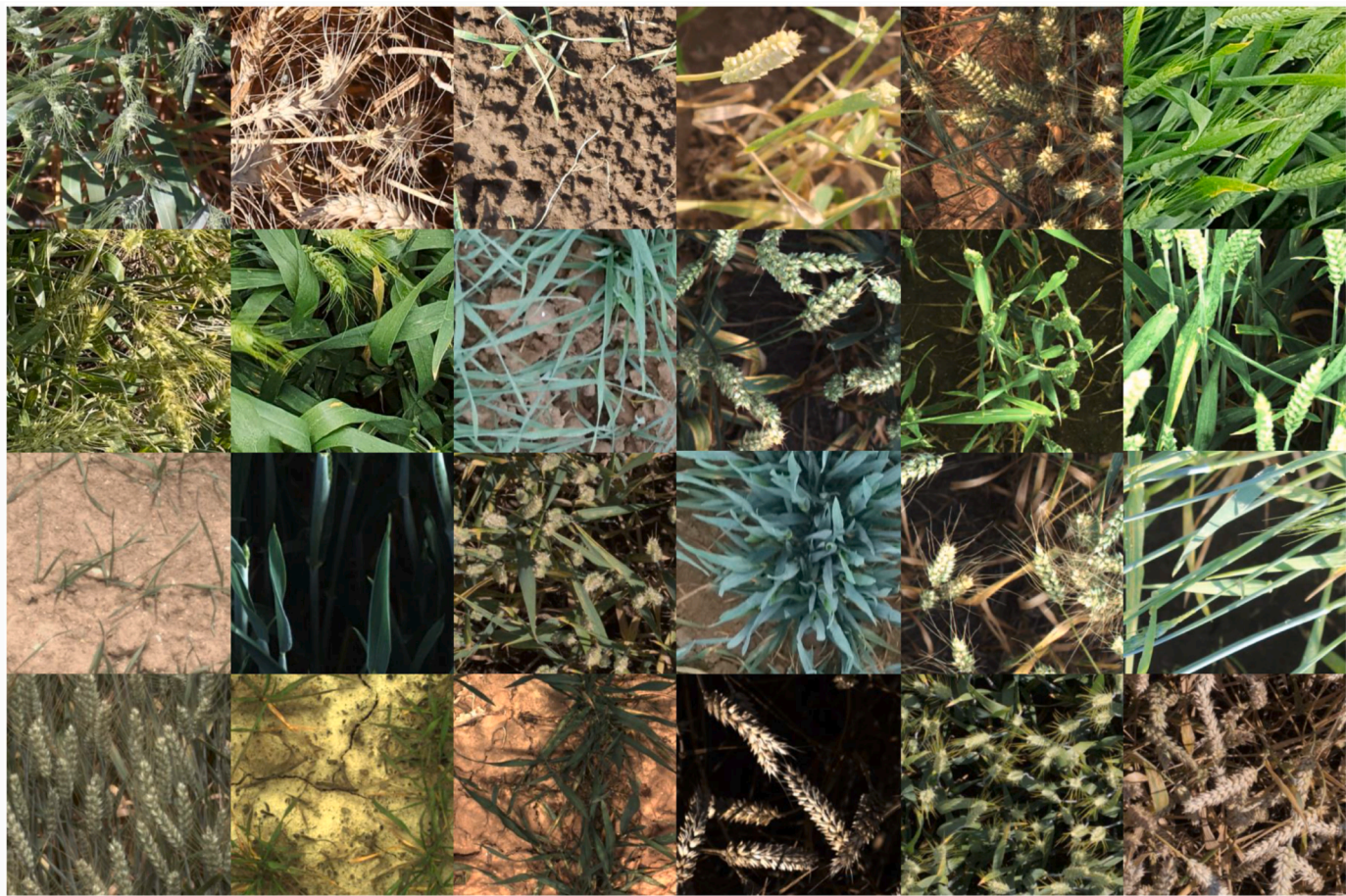


Fig. 2. Representative samples from our proposed GWFS dataset.

Table 3
Imaging tagging overview.

Tag Name	Description	Possible Value
Institute	The institute that contributes this image.	INRAE, ETHZ, USASK, Arvalis, RRES, NJAU, UQ, ULiege, Utokyo, CIMMYT.
Name	The unique image name.	Not Applicable.
Size	The width and height of the original image.	(4800 × 3200), (2048 × 2048), (4096 × 3000), (3456 × 5184), (5184 × 3456), (5634 × 3753), (2040 × 2044), (4080 × 2704), (4080 × 3200), (3296 × 2472), (2560 × 2048), (1024 × 1024), (5184 × 3456).
Anthers	The existence of anthers.	True, False.
Bending	The existence of ear bending.	True, False.
Lighting Stage	The existence of shadow. Phenological Stage.	True, False. Emergence, Vegetative, Stem Elongation, Ear Emergence, Early Filling, Early Senescence, Late Senescence, Maturity.

sheath requires either its edge or the attached leaf blade is visible (indicated by “S” in Fig. 3 d). In many cases, it was not possible to decide if the structure was sheath, bare stem or peduncle. The wheat head label specifically encompassed only spikelets, excluding awns. Awns are distal bristle-like extensions of the lemma surrounding the florets of wheat [79]. Thus, while spikelets including glumes and florets were labelled as part of the wheat head, awns were treated as “invisible” features (Fig. 3 f, h). There are practical challenges in annotating individual awns in all imaging scenarios, particularly when awns appear blurry or lack distinct contours. Thus, when awns overlapped with the targeted organ, they

were treated as if they were absent and the polygon was drawn across the organ in the background (Fig. 3 f, h). Detached senescent plant material, such as debris or residue resulting from no-till practices, was excluded from annotation and classified as background. All targeted entities were labelled as long as they were still attached to the wheat plant regardless of their colour, i.e. chlorotic or necrotic tissue was also labelled.

The annotation process was exclusively focused on wheat, with weeds deliberately left unannotated and classified as background. All other non-wheat objects were similarly disregarded and annotated as background (Fig. 3 b, h).

2.5. Baseline segmentation models development

DeeplabV3Plus [80] is a classic semantic segmentation framework based on convolutional neural networks (CNN) that employs an Encoder-Decoder architecture. The model builds upon the strengths of Deeplabv3, which leverages Atrous Convolution to explicitly control the resolution of feature maps and adjust the reception field. In Deeplabv3Plus, encoder features are first up-sampled bilinearly by a factor of 4 and then concatenated with the corresponding low-level features from the backbone network. 1 × 1 convolutions are applied to low-level features, reducing the number of channels to reweight rich contextual encoder features and simplify training. After the concatenation, the model refines these combined features using a series of 3 × 3 convolutions, ensuring the integration of detailed spatial information and high-level semantic context. The decoder finalises the segmentation mask with a simple bilinear upsampling operation by a factor of 4, delivering high-resolution predictions. This seamless combination of

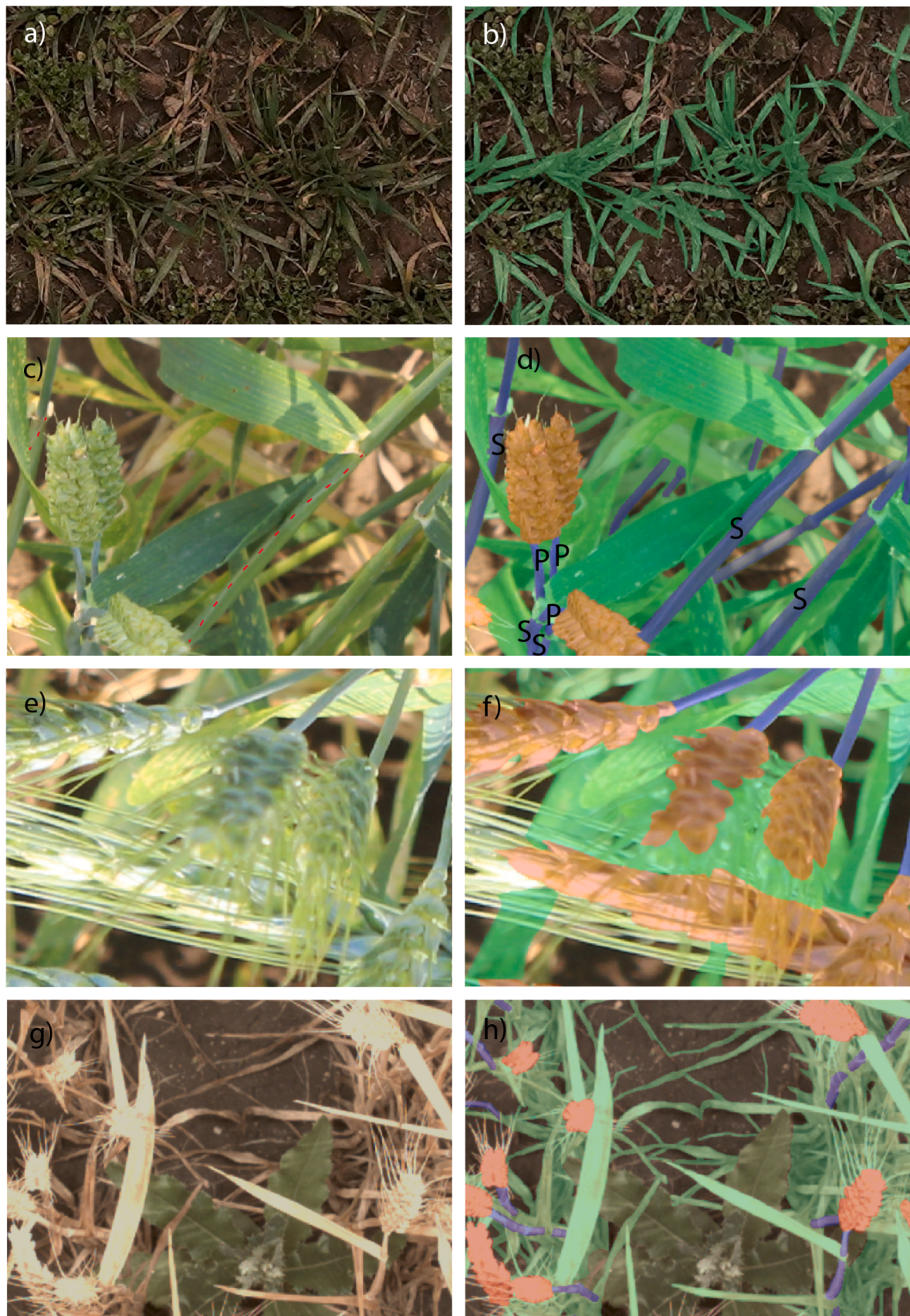


Fig. 3. Examples for the labelling process with original images (a, c, e, g) and labelled images (b, d, f, h) showing the segmentation masks for leaves (green), stems (purple) and spikes (orange). Red dotted lines in a) show the edge of the leaf sheath wrapped around the stem. A leaf sheath (S) as part of the stem can be either recognized by this edge or its connection to a leaf blade. A peduncle (P) is a part of the stem located between a visible spike and the collar of the flag leaf. For our analysis awns were treated as invisible by drawing the segmentation masks above them (f, h). Weeds and crop residues were considered background (b, h).

multiscale context aggregation through the encoder and spatial detail recovery in the decoder positions Deeplabv3Plus as a robust and flexible solution for semantic segmentation tasks.

Segformer [63] processes an input image of size $H \times W \times 3$ by first dividing it into 4×4 overlapping patches. These patches are fed into a hierarchical Transformer encoder to extract multilevel features, leveraging an Overlapped Patch Merging strategy to ensure spatial continuity and capture richer local context. The encoder utilises efficient self-attention mechanisms by applying dimension reduction, which reduces the time complexity of the self-attention mechanism. This significantly optimises computational efficiency with minimal impact on segmentation performance.

The extracted features are then processed by a lightweight multi-layer perceptron decoder. Unlike traditional hand-crafted designs such as Deeplabv3Plus, this approach simplifies the decoding process while improving the effective receptive field, enabling precise and efficient segmentation. To cater to diverse performance and resource requirements, Segformer introduces a family of Mix Transformer encoders (MiT-B0 to MiT-B5), which share the same architecture but vary in size, offering flexibility in balancing computational cost and segmentation accuracy. By integrating innovative design choices with practical adaptability, Segformer delivers a robust and efficient solution for semantic segmentation tasks.

2.5.1. Impact of distribution shift on segmentation performance

We conduct experiments under two different data-splitting settings. (1) **Random Split**: In this setting, we randomly split the data into a training set (70 %), a validation set (10 %) and a test set (20 %). (2) **Region Split**: In this setting, we utilised data from Arvalis, CIMMYT, ETHZ, INRAE, NJAU, RRES, and ULiege CRA-W as the training set, data from UTokyo as the validation set, and data from UQ as the test set. The UQ test set is rather challenging due to the massive diversity in genotypes and Australian growing environments and imaging conditions. For both of the settings, the validation mIOU was used to select the best checkpoint, which was then used for testing.

2.5.2. Impact of training data scale and model size on segmentation performance

We conducted two sets of experiments to investigate how the size of the training dataset and the number of model parameters influence the performance of the segmentation. To investigate the relationship

between the size of the training dataset and the performance of the model, we trained SegFormer-B0 using progressively larger subsets of the full training dataset. Specifically, we sampled subsets containing 1 %, 5 %, 10 %, 20 %, 30 %, ..., 100 %, of images from the full training set to train the model. To assess the impact of model size on segmentation performance, we trained SegFormer-B0 to B5 using the full training dataset. As the model progresses from B0 to B5, both the number of parameters and computational cost increase, allowing us to analyse how model complexity affects segmentation accuracy.

2.6. Evaluation metrics

Mean Intersection over Union (mIoU) measures the level of overlap between the predicted mask and the ground truth mask. Specifically, we have:

$$\text{mIoU} = \frac{1}{C} \sum_{c=1}^C \frac{|P_c \cap G_c|}{|P_c \cup G_c|}, \quad (1)$$

where P_c and G_c denote the predicted mask and ground truth mask of the c -th class.

Mean Pixel Accuracy (mAcc) focuses on the pixel-wise accuracy for each class, which can be defined as:

$$\text{mAcc} = \frac{1}{C} \sum_{c=1}^C \frac{|P_c = G_c|}{|G_c|}. \quad (2)$$

3. Results

3.1. UMAP visualisation of image diversity

The diversity of the training dataset is critical for the generalisation capacity of segmentation models. The diversity of a crop dataset can be affected by multiple factors, such as differences in phenological stages, lighting conditions, background, growing environment conditions, and genotype. In this work, the distribution of GWFS images was analysed using the UMAP technique on image features extracted by an ImageNet-pretrained ResNet-50 model. The image embeddings of the top two UMAP components visualise the distribution of all labelled GWFS data. The visualisation of images in the latent space was either colored by the phenological stage (Fig. 4 a) or institution (Fig. 4 b). For the

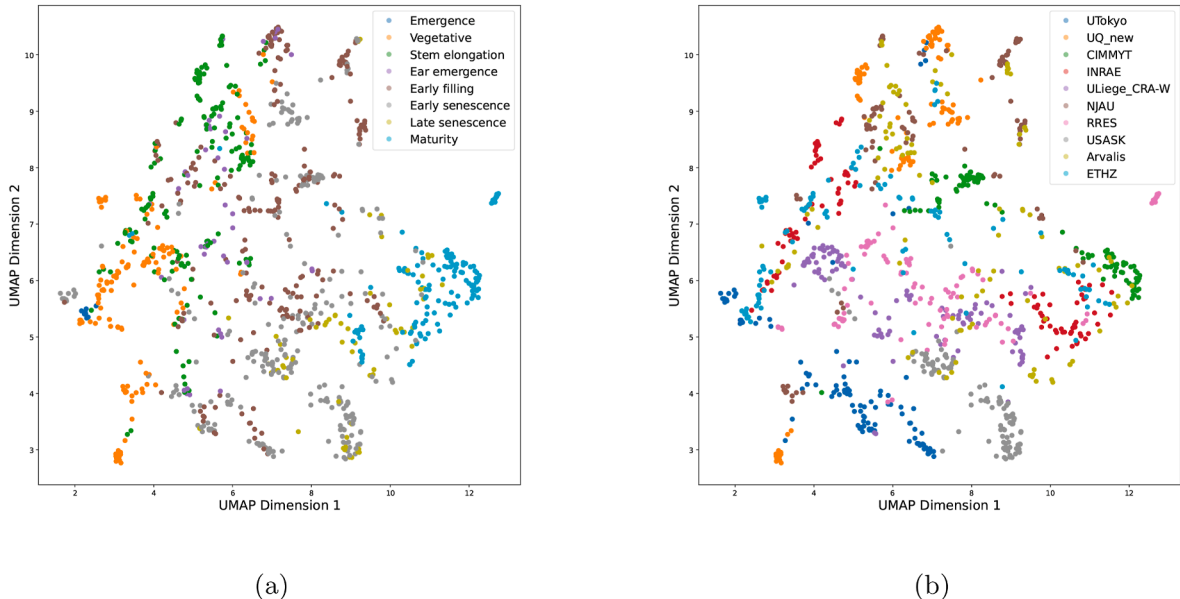


Fig. 4. The UMAP visualisation of GWFS labelled images colored by (a) phenological stage and (b) institution.

Table 4
Glossary of relevant terms used.

Term	Abbreviation	Description	Reference
Labels		GWFS context: polygons used to generate masks for the different organs.	
Tags		GWFS context: keywords added to an image, such as the developmental stage displayed.	
Tiles		GWFS context: subsamples from the original image that are used for labelling and training. Standard tiles are squares, often with 256×256 pixels.	
Uniform Manifold Approximation and Projection	UMAP	Dimension reduction technique that can be used for latent feature visualisation	McInnes, L, Healy, J, UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction, ArXiv e-prints 1802.03426, 2018
Intersection over Union	IoU	IoU is a metric of segmentation performance of the area manually labelled (A) vs. the area detected by the model (B). The intersection size between the two is divided by their union size. intersection union	https://en.wikipedia.org/wiki/Jaccard_index
Proximal sensing		Sensing from proximity. Sensors are typically operated hand-held or mounted on poles, gantries, or ground-based vehicles.	
Imaging setup		The combination of vector (carrier systems), sensor and lens (i.e. RGB camera-lens combination) as well as the working distance and camera angle defining the field of view.	https://en.wikipedia.org/wiki/Field_of_view vor definition of the field of view.
Phenology		The study of periodic events in biological life cycles and how these are influenced by seasonal and interannual variations in climate, as well as habitat factors (such as elevation)	https://www.merriam-webster.com/dictionary/phenology
Canopy		Branches (stems), leaves, and spikes (inflorescences) of a population of plants growing on a piece of land.	Adapted from different sources ^a
Crop ontology	CO	Provides descriptions of agronomic, morphological, physiological, quality, and stress traits along with their definitions and relationships.	https://croponontology.org
Leaf area index	LAI	Trait, which characterise plant canopies, are defined as leaf green area per unit of surface area	CO_321:0000184
Initiation of booting	Boot	Phenological period prior to spike emergence where the flag leaf is fully developed	CO_321:0000191
Anthesis	Ant	Phenological period when pollination occurs in wheat	CO_321:0000121
Heading	Hd	Phenological period from the time of emergence of the spike tip from the flag leaf until the spike has fully emerged	CO_321:0000007
Maturity	Mat	Phenological stage when wheat stops remobilising assimilates to the spike	CO_321:0000022

^a We combined definitions from <https://dictionary.cambridge.org/dictionary/english/canopy> and <https://www.ebi.ac.uk/ols4/search?q=canopy> as no matching definition was found for the case of crop canopies.

phenological stages, the first UMAP dimension shows a clear clustering in the sequence of stage progression. With the first UMAP feature increased from 0 to 13 (Fig. 4 a, dimension 1), the growth stage generally progresses from emergence and vegetation towards senescence and maturity. Most of the images taken from emergence to stem elongation clustered at a value below 7, while the other extreme images around late senescence and maturity clustered above 7. Images containing ears, i.e., starting from ear emergence, showed values above 4 in dimension 1. The clustering in the second dimension tended to be driven by the institution providing the images. Images from USASK, ULiege, and RRES generally have the second UMAP feature valued below 8, while most of the UQ and NJAU features are above 8. By viewing the images of these institutions, this phenomenon could be attributed to the difference in the lighting conditions. The datasets of UTokyo and NJAU cover a wider range of dimension 2, while the dataset of Arvalis is the only one spanning almost the entire latent space.

3.2. Balance of developmental stages and labelled classes

Images were tagged with the approximate developmental stage estimated from the image (i.e. stages were not recorded as ground truth in the field). Although the aim was to balance all the phenological stages (for a definition see Table 4), this was not possible for all datasets. An analysis of the tags across the whole dataset revealed an uneven distribution across different phenological phases (Fig. 5, a). Notably, 'early filling' and 'early senescence' were the most frequently observed (224 and 282 images, respectively), while 'emergence' was the least represented stage (8 images).

We also evaluated the balance of the labelled classes. At the image level, background (BG) and leaves were the most prevalent, appearing in 1090 and 1080 images, respectively (Fig. 5, b). In contrast, the stems

and heads were present in 847 and 739 images, respectively, since they predominantly emerge in the later phenological stages. At the pixel level, a more pronounced class imbalance was evident, with the leaves occupying the largest proportion of pixels, followed by the background, while the stems and heads account for significantly fewer pixels (Fig. 5, c). Thus, although heads and stems appear frequently at the image level, they still constitute only a small fraction of the total pixel distribution, compared to leaves and background.

3.3. Baseline segmentation models

Concerning the sampling strategies to split images into training, validation and test sets, there was a noticeable performance gap between the random split vs region split strategy (Table 5). We attribute this discrepancy to the distribution shift between the training and test data in the region split setting. Among the evaluated models, Segformer consistently outperformed DeepLabV3plus, achieving a 2.8 % higher mIoU in the random split setting and an 8.5 % improvement in the region split setting. Notably, in the Region split strategy, Segformer performed substantially better for head and stem classification than did DeepLabV3plus with little difference in estimation of background and leaf. These results highlight the superior capability of Segformer in addressing the wheat organ segmentation task, especially under distribution shifts.

Concerning the size of the training data, the overall model performance sharply increased between 1 % and 60 % (i.e. 460 images) of the training data and plateaued thereafter (Fig. 6, a). Above these 460 images, there were only marginal improvements as the dataset size approached 100 %. The heads and leaves were well segmented with only 10 % of the data ($IoU > 75 \%$) and the model performance only progressed slowly when more data were used. The segmentation of stems

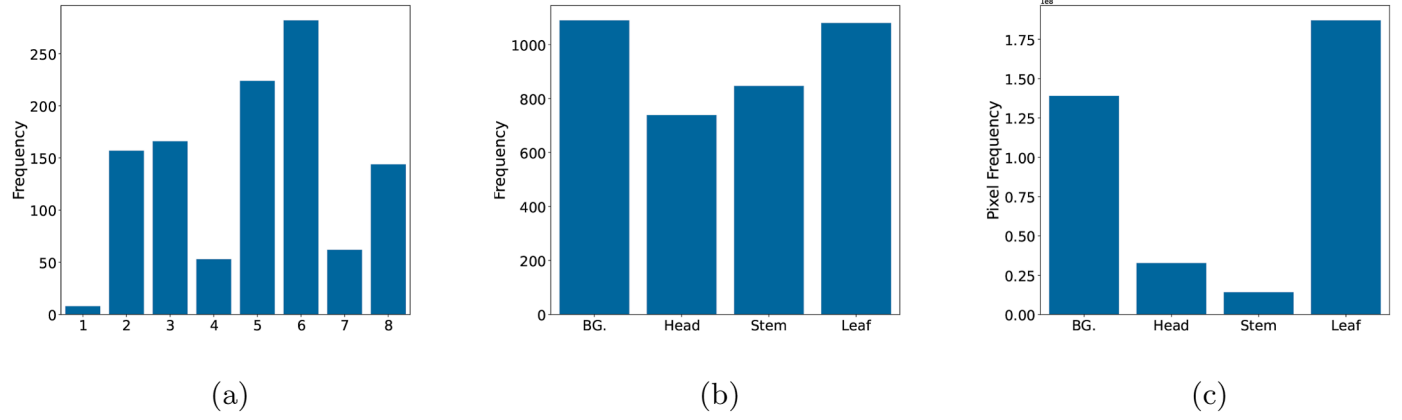
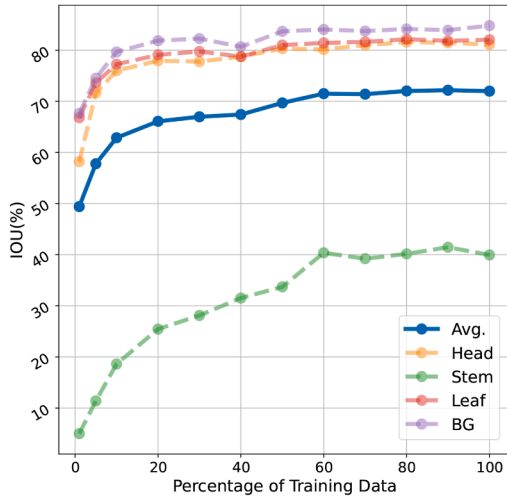


Fig. 5. Statistics of the GWFS dataset: (a) The distribution of growth stages for each image, where 1 – 8 on the x-axis indicates ‘Emergence’, ‘Vegetative’, ‘Stem elongation’, ‘Ear emergence’, ‘Early filling’, ‘Early senescence’, ‘Late senescence’, and ‘Maturity’, respectively. (b) The class occurrence at the image level. (c) The class occurrence at the pixel level.

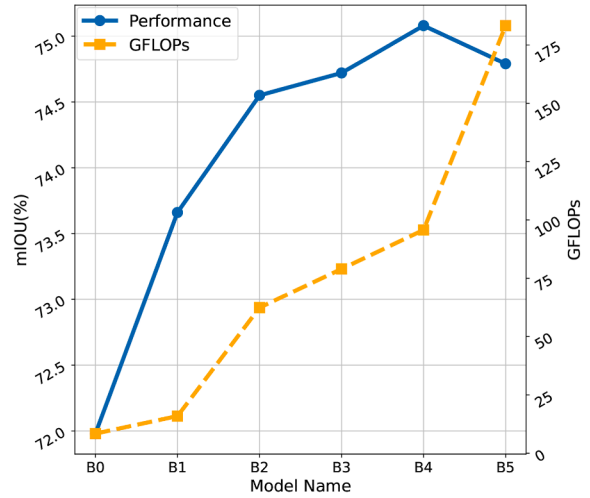
Table 5

Comparison of Mean Intersection over Union (mIoU) and mean accuracy (mAcc) metrics for Deeplabv3plus (R101) and Segformer (B1) across Random- and Region-based data splittings.

	Segformer (B1)				DeeplabV3+ (R101)			
	Random		Region		Random		Region	
	IoU	Acc	IoU	Acc	IoU	Acc	IoU	Acc
Background	84.51	92.75	75.46	88.71	81.59	89.88	74.32	82.11
Head	82.85	90.14	66.11	86.84	81.46	89.05	46.25	48.27
Stem	44.92	53.89	19.23	20.95	39.40	46.73	6.64	7.06
Leaf	82.35	90.44	81.75	89.36	80.93	90.81	81.4	95.14
Average	73.66	81.81	60.64	71.47	70.85	79.12	52.15	58.15



(a) Training Data Scale vs. Model Performance



(b) Model Scale vs. Model Performance

Fig. 6. Effect of (a) training data size and (b) segmentation model size on Segformer model performance. The effect of training data size on the intersection over union (IoU) was evaluated for all object classes including the background (BG). The segmentation model size performance was evaluated as mean IoU (mIoU) and Giga Floating Point Operations per Second (GFLOPS). The higher GFLOPS indicates heavier computational complexity. Here B0 indicates the smallest Segformer model while B5 indicates the largest Segformer model.

improved with increases to 60 % of the training data but the performance plateaued at a low IoU of 40 %.

To examine the impact of model size, we trained SegFormer models of varying capacities (B0 to B5) on the full training dataset and analysed their performance (Fig. 6, b). The Giga Floating Point Operations per Second (GFLOPS) served as an indicator of computational complexity,

where higher values denote greater computational demands. The results indicate a steady improvement in performance from B0 to B4. However, when training SegFormer-B5 on the full dataset, performance degradation was observed. We attribute this decline to the insufficiency of training data to adequately support the significantly larger parameter space of SegFormer-B5, which nearly doubles that of B4.

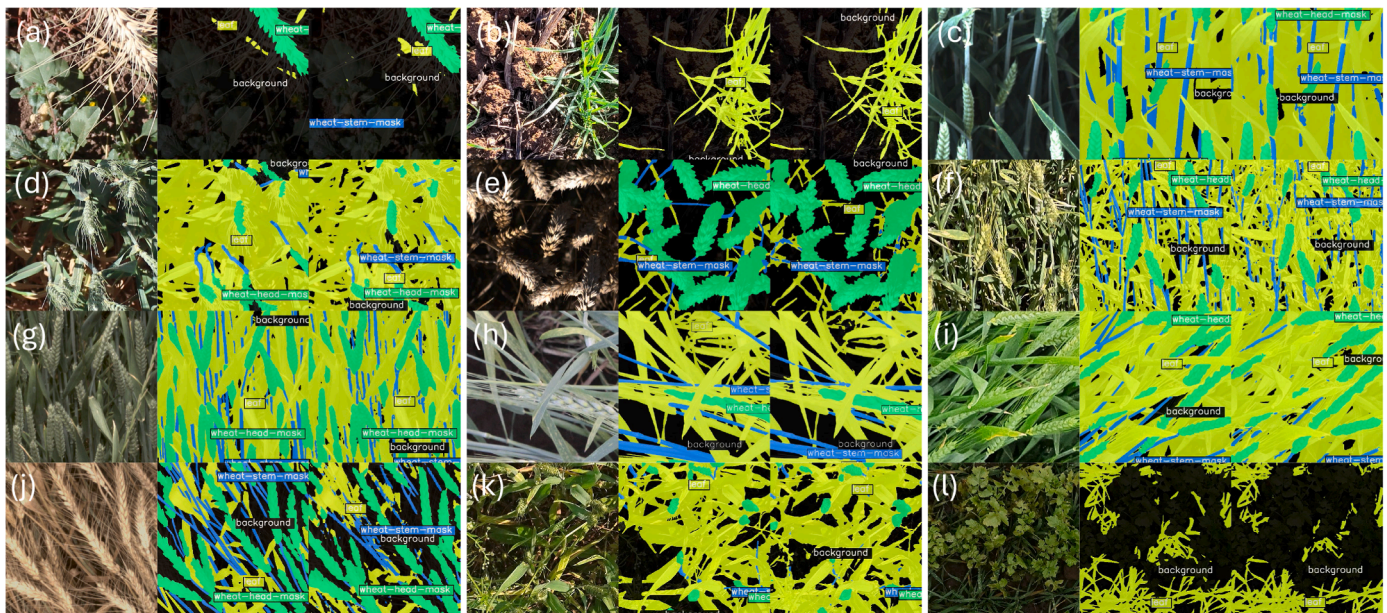


Fig. 7. Visualisation of original image (left strips), ground truth (centre strips) and prediction results (right strips) from Segformer-b1.

3.4. Visual inspection of segmented images

Key requirements of the model were to detect organs independently of their colour and distinguish them from weeds or plant residues on the ground. We conducted a systematic review of all five images that contain weeds in the test set (Random Split, as described in Section 2.5.1). The weeds were of different types with different leaf shapes and were mostly classified as background (Annex, Fig. S2). A sample is displayed in Fig. 7 a, k and l. Late in the season, weeds frequently germinate between rows and obscure the senescence signal. Similarly, the separation of plant residues from wheat plants (Fig. 7, b) is an important advantage. When awns are present, the segmentation of the spike without the awns (Fig. 7 a, d, f, h, i, j, k) is a useful feature that can assist in the approximation of spike volumes. To evaluate the segmentation of senescing canopies, we selected images from early senescence to maturity within the random split test set (as detailed in Section 2.5.1). This set representing later stages comprises 73 images. In general, the model maintains strong performance (Table S2). However, performance declined relative to the full test for the background (IoU of 73.5 vs. 84.5 for late stage vs. full set) and leaves (IoU of 73.7 vs. 82.3 for late stage vs. full set). This is likely due to the increased visual similarity between senescing leaves and background elements (e.g., soil or dried residue), as well as reduced structural distinctiveness in ageing foliage, making accurate segmentation more challenging.

4. Discussion

4.1. The challenges of organ labelling in complex canopies

We collected 52,078 images from 67 different field sites worldwide with a ground sampling distance (GSD) between 0.09 and 0.71 mm per pixel. The UMAP visualisation confirms the need for such diversity, as it shows clustering by institution and developmental stage.

The collaborative effort to sample this diverse set of images and to design the labelling strategy was essential to the success of the work. All participating institutions operate imaging setups collecting images of wheat canopies in the sub-millimeter range and are at the forefront of enhancing the in-depth analysis of complex canopies. A first step was the decision of which canopy features could be targeted, given the available spatial resolution. In wheat, ground sampling distances below 10 mm permit good estimates of canopy cover and leaf area index or crop

density, while GSDs below one mm are needed to detect individual leaves of emerging seedlings [81]. The given sub-millimeter resolution was sufficient to label the targeted organs. Reliable labelling of awns would likely require GSDs below 0.2 mm and substantially more labour for manual labelling. New sensors with higher resolution will allow for awn segmentation, even when operating with the same carrier system.

Given the experience in GWFSS, the term “plant organ” needs to be understood in the context of imaging constraints. With this regard, we would like to add a “sidenote” to protocols for the minimal requirements to describe plant phenotyping experiments. The MIAPE 1.1 [82], release states that “Observed variables, traits, methods and scales are each identified by name, and may have a reference to the corresponding ontology concept (ideally from the Crop Ontology)”. However, in the Crop Ontology database,¹² entities upon which traits are measured are not always indexed with their own identifiers. For example, while stem colour (CO_321:0000973) is defined as “colouration of the stems,” the entity “stem” itself is not defined. GWFSS aims to extract organs as entities on which traits will be measured. For this reason, we prefer the BRENDA Tissue Ontology (BTO)¹³, which is a vocabulary for the source tissues. As such, it focuses on a detailed description of the entity rather than the trait. The segmentation of the entities is the first step towards deriving new phenotypes. Moreover, our digitally extracted “organs” do not quite comply with classic ontology terms: heads exclude their spikes, and stems include leaf sheaths. One challenge is therefore how to define a trait which is based on more than one organ part (e.g. stem = true stem plus leaf sheath), given that trait ontologies are typically structured hierarchically. Along these lines, Celestina et al. [83] have identified the need to reconsider the classical growth scales, such as Zadoks [84] or the derived unified BBCH scale [85] to fit the needs of image-based phenotyping. They draw up a list of the development stages that need to be assessed by destructive sampling and the stages that can be assessed in a non-destructive manner. We believe that their phases of the Population of Culms Development Scale (PCDS) may be enhanced by image-derived phases. At least heading and physiological maturity can be digitally measured based on models derived from GWHD and GWFSS as we will discuss below.

¹² <https://cropontology.org/>.

¹³ <https://bioportal.bioontology.org/ontologies/BTO>.

Table 6
Summary of wheat segmentation datasets.

Dataset	Size	Resolution	Focus
GWFS	52,078 images	Various	All organs, all growth stages, geographic diversity, genotype variation
GWHD [16]	6510 images	Various	Wheat Head detection
WESS-Dataset [65]	120 images (6500 tiles)	4592 × 3448 px (256 × 256 tiles)	Wheat heads
EarSegNet [66]	160 images	5184 × 3456 px (2500 × 2500 tiles)	Flowering wheat heads
Najafian et al. [36]	Limited manual + synthetic	12/48 megapixel	Wheat heads
Deng et al. [72]	370 images (25,530 tiles)	3000 × 2000 px (256 × 256 tiles)	Leaves
Anderegg et al. [73]	206 tiles	2400 × 2400 px	Stem elongation
Liu et al. [71]	2200 images	Not specified	Fusarium Head Blight
Niu et al. [89]	450 images	3472 × 3472 px	Head damage

4.2. GWFS compared to other datasets for semantic segmentation of wheat organs

The availability of extensive open datasets has been crucial to remarkable progress in applications of modern computer vision in agriculture. Various datasets have been released to facilitate computer vision applications, such as for fruit detection [86], weed management [87], green coverage estimation [9], and plant disease identification [88]. For semantic segmentation, pixel-level annotation remains a cornerstone of segmentation tasks, but it is notoriously labour-intensive and expensive. Intensive collaboration among public and private institutions is needed to generate sufficiently large, diverse, and consistently labelled data.

The novelty of the GWFS dataset is its diversity, its coverage of all stages, and labelling of all organs of wheat. The GWFS dataset contains fewer annotated images than GWHD [16], which labelled 6510 images from 16 institutions. However, GWFS dataset samples the whole growing season while GWHD focuses on head detection during flowering, grain filling, and ripening. Moreover, for the GWFS dataset, we decided to supply a large set of 52,078 images to enable users to pose solutions to other questions by applying their labelling. A similar dataset assembled by the global-rice dataset consortium is underway for rice [73].

For semantic segmentation, there are several smaller datasets available which are not included in GWFS, mainly because they focused only on specific organs of wheat. Although most studies trained segmentation models, we will summarise only the characteristics of the annotated datasets as the most valuable part of the studies (see Table 6). With regard to wheat heads, several smaller datasets were collected in addition to the GWHD dataset [36,65,66]. Few researchers have collected datasets to train leaf segmentation during the early [72] or late phases of development [73]. The latest trend is to target features within organs within complex canopies in the field, such as spikelet segmentation [89], Fusarium head blight [71], or leaf diseases [90]. For such annotations of wheat head damage, a concise and coordinated labelling of diverse datasets may be highly valuable.

4.3. Organ segmentation - a path to integrative traits

Field phenotyping is often considered in the context of spatial and temporal scales. We believe that the ability to track organ development throughout the growing season will set a new standard for phenotyping. It will aid breeders, variety testers, or researchers in evaluating genotypes for improved canopy architecture, source-sink balance and resistance to environmental and disease stressors.

The ability to follow the different wheat organs through the season enables the testing of more complex phenotypes. For example, our model would permit quantification of the changes in tissue reflectance

as they change from green to chlorotic to necrotic. In wheat, this was done, for example, using shallow learners based on colour spaces using a support vector machine classifier [19] or a multiclass random forest classifier [73]. Canopy segmentation followed by reflectance indices or the mentioned classifiers opens new possibilities for field phenotyping. For example, changes in leaf colour can be used to quantify canopy damage in winter due to frost or diseases. Later, it may be used to quantify nitrogen status.

Stem (or more precisely peduncle) senescence is a measure of physiological maturity of wheat [91]. It is time-consuming to assess and, therefore, rarely reported. Using a small semantic segmentation training dataset, Anderegg et al. [43] tracked the senescence process of leaves, stems, and heads through grain filling and showed that stems were the last organ to be senescent. In this work, the genotypes differed in their dynamics and timing of leaf, head, and stem senescence, highlighting the relevance of tracking the organs separately. The GWFS training dataset is a major step forward in measuring physiological senescence and separating different canopy senescence processes.

We acknowledge that proximal sensing is currently not the primary choice in cases where thousands of plots are to be evaluated and breeding programs are more likely to require high-throughput remote sensing by means of unmanned aerial vehicles (UAV). Remote sensing by UAV equipped with lower pixel resolution multi-spectral sensors allows the estimation of traits like LAI and canopy senescence. However, these sensors are not ideal to study the different organs in the canopy. Plant organs may have substantially different proportions in the canopy, influencing its reflectance spectrum. Variation with time occurs due to environmental effects on crops (leaf rolling, frost, heat, drought, nutrition, and pest effects), stem elongation, spike appearance, and bending or lodging of spikes during grain filling. The transferability of reflectance spectra from one season or trial to a different one may be influenced by environmental conditions or diseases. Although proximal sensing has lower throughput, it can complement aerial measurements with information from the organ scale. This might enable upscaling from proximal to remote sensing. Alternatively, high-resolution RGB imaging is becoming readily available on UAV with cameras with a resolution ranging from 60 to 120 Mpx, though UAV are slow and difficult to localise when flown close to plots to take still shots. A solution is to fly UAVs closer to the plots and use video at high-shutter speeds to avoid blurring the image or disrupting the canopy with propeller downwash [92] were able to achieve a GSD of 0.13m flying a 20 Mpx camera using video at ca. 5m above the canopy at a speed of 2–3 m/s. Finally, there is an increasing number of phenotyping robots available, as well as higher resolution RGB UAV cameras that allow organ phenotyping of larger numbers of plots.

We recognise that RGB datasets are only a piece of the larger phenotyping toolbox: Other technologies could and should be used in combination with data fusion methods. For example, LiDAR has been

used to estimate biomass and crop growth rate [93]. Future RGB datasets collected at the same time as LiDAR point clouds can be used to study biomass accumulation at deeper hierarchical levels, such as biomass partitioning in the different organs. Furthermore, RGB combined with thermal imagery can be used to assess abiotic stresses [94]. Breeders could focus on targeted breeding for specific organs and select new genotypes with increased water use efficiency, to exemplify a few cases where RGB organ segmentation can be made the most of.

4.4. Baseline model to guide the size of the training data

The main contribution of this work relates to the creation of a large open-access database of images capturing diverse field-grown wheat plots. However, without a basic model it was difficult to judge how large the training data should be. According to the trained baseline model, the 1096 labelled images are sufficient for a segmentation of leaves and heads. For any organ, including stems, the increase in model performance levelled out when more than 60 % (ca. 600 images) of the totally available data were used for training. This indicates that either a massively larger amount of data or a different labelling or training strategy might be needed for substantial further improvement. Model performance was reduced when performing a region-specific data split.

Among the two evaluated models, we chose Segformer. It achieved good segmentation for leaves, heads and background (IoU 80 %), but the segmentation of the stems underperformed (IoU = 44.92 %). The low performance of the stem segmentation may be attributed to various reasons: stems are thin, partially occluded by leaves, and a limited proportion of total pixels compared to the other organs. During booting and spike emergence (i.e. the expansion of the growing spike within leaf sheaths), the cylindrical structure of the stem is lost and the foil-like structure of the sheaths becomes obvious. This may lead to low model performance during this stage, particularly for stems. Furthermore, annotation difficulties related to distinguishing stems from some soil residues and senescent or rolled leaves may produce label noise. Improving stem segmentation may require exploring loss functions associated with class imbalance or incorporating massively more images with annotated stem masks.

These findings highlight a critical requirement: current models for segmenting organs in plants are highly dependent on extensive, high-quality annotations to achieve strong performance. To address this, research should prioritise the development of annotation-efficient solutions that maintain high performance with fewer labelled samples. Promising directions include leveraging self-supervised learning, which utilises unlabelled data and semi-supervised or active learning, which strategically selects a minimal number of samples for annotation while maximising model learning.

Without being quantitative, visual observation of the segmented image confirmed that wheat plants were segmented independently of their colour while weeds or plant residues were classified as background. This is a great step forward. Earlier models segmented plant tissue to a large degree based on colour. For example, in the case of Eschikon wheat segmentation training (EWS) [18], necrotic leaves were no longer detected as part of the canopy. This led to a decrease in canopy cover after winter due to necrotic leaves that suffered freezing damage [18]. Although the VegAnn model [9] included necrotic parts of plants in its segmentation while excluding crop residues, it was trained on a large number of different species and did not exclude weeds [34].

4.5. Value of the full GWFSS dataset without labels

Self-supervised learning (SSL) methods can leverage a large amount of unlabelled image data as a pre-training procedure to better initialize or condition a deep learning model for a subsequent downstream analysis task. Ogidi et al. [34] found that a diverse source dataset in the same domain or similar as the target dataset combined with SSL can maximize performance in downstream plant phenotyping tasks. Our large, unlabelled dataset of 52,078 RGB images is meant as a training set for SSL methods. The idea of this dataset is to provide it as set for SSL while using the labelled data for validation and testing. The dataset can be also used for further labelling without the need to collect the data. It may be used for further segmentation tasks such as spikelet segmentations or other canopy features.

Beyond organ segmentation, the extensive dataset documented here provides a dynamic platform to develop predictive models that can capture temporal and spatial variability across multiple years and environments. By integrating environmental data with image-derived traits and machine learning methods, such as random forest regression or XGBoost, which can handle a vast array of predictor variables, breeders can target more complex traits such as radiation use efficiency, harvest index and yield. More fundamental underlying traits can be capable of better account for genotype-by-environment interactions and permit breeding programs to optimise their pipelines on a global scale.

4.6. Conclusion and outlook

Segmentation models will reduce the subjectivity of field observations by leveraging the generation of large and consistent datasets. Organ segmentation will enable the extraction of a range of additional traits from complex canopies. Such information is needed to advance our understanding of the interaction of *Genotypes* with the *Environment* they grow in and the *Management* practices they receive (often abbreviated as GxExM). To achieve robustness, the training data for the segmentation models needs to be large and diverse. Models derived from the GWFSS dataset will likely outperform many models derived from labelling in single experiments. As a further step, datasets from other small-grain cereals, such as barley, could be considered to enhance the training data.

While organ proportions will be directly available from the GWFSS-derived segmentation models, other traits will need to be developed and calibrated through secondary processing. For example, sensor fusion may allow one to integrate organ information derived from point sensors or sensors with lower resolution. Moreover, organ information may complement canopy-level traits derived from remote sensing.

Many of our images were derived from fully automated platforms, and these installations are often prototypes. With the advancement of agricultural robots, such platforms will become affordable for a greater community in the near future. But also smartphones or specifically designed hand-held devices, will bring image-based phenotyping to a greater community. This will leverage new possibilities for common research projects in science and citizen science communities.

Author contributions

The GWFSS consortium consisted of different working groups that focused on conceptualisation and steering, data collection and data supply, labelling, data curation, training the base model, and writing (Table 7). The display of the author's contribution was inspired by.¹⁴

¹⁴ <https://www.nature.com/nature-index/news/researchers-embracing-visual-tools-contribution-matrix-give-fair-credit-authors-scientific-papers>.

Table 7

Author contribution.

First Name	Last Name	Steering	Data supply	Data curation	Labeling	Base model	Writing	Funding
Benoît	De Solan	hi	mid		hi		mid	hi
Lucas	Bernigaud Samatan		mid		mid		lo	
Safaa	Ouahid		mid	hi	hi		lo	lo
Andrea	Visioni	mid	hi	hi	mid		mid	
Carlos A.	Robles-Zazueta	mid	mid				hi	
Matthew P.	Reynolds	mid					lo	lo
Francisco	Pinto	mid	mid				mid	
Ivan	Perez-Olivera		mid				lo	
Shouyang	Liu	lo	mid				lo	lo
Chen	Zhu		mid				mid	
Marie-Pia	D'argaignon		mid				lo	
Raul	Lopez-Lozano	lo					mid	lo
Marie	Weiss	lo					mid	
Radek	Zenkl		hi	hi	hi		hi	
Andreas	Hund	hi			mid		hi	hi
Norbert	Kirchgessner		lo				lo	
Afef	Marzougui		mid				lo	
Lukas	Roth		lo				mid	
Alexis	Carlier		lo				lo	
Sebastien	Dandrifosse		hi				hi	
Benjamin	Dumont		mid				lo	lo
Benoît	Mercatoris	mid					lo	lo
Zijian	Wang	hi	hi	hi		hi	hi	
Scott	Chapman	hi	hi		lo		mid	hi
Javier	Fernandez		mid				lo	
Zhi	Chen		mid				lo	
Ian	Stavness	hi		hi			hi	lo
Keyhan	Najafian		mid				low	
Wei	Guo	hi	mid	hi		mid	mid	lo
Haozhou	Wang		mid				lo	
Latifa	Greche		hi				hi	
Malcolm	Hawkesford	hi					lo	lo
Nicolas	Virlet		hi				hi	
Etienne	David	mid					mid	
Alexis	Comar	hi			hi		lo	
Joss	Gillet	mid		hi			mid	
Kamran	Irfan	mid					lo	

Funding

Global wheat was directly supported by Analytics for the Australian Grains Industry (AAGI). AAGI (UOQ2301-0100PX) is a Strategic Partnership between the Grains Research and Development Corporation (GRDC), Curtin University, The University of Queensland and the University of Adelaide. Other project and associate partners also support the initiative; Arvalis, France; Phenet (European Commission [95]); EMPHASIS-GO (European Commission [96]); Delley Seeds and Plants Ltd, Switzerland; Deutsche Saatveredelung AG, Germany; ASUR plant breeding, France; Plant Phenomics journal, China.

Funding of individual projects of partners

CIMMYT: The International Wheat Yield Partnership (IWYP); the Heat and Drought Wheat Improvement Consortium (HeDWIC); the Accelerating Genetic Gains in Maize and Wheat (AGG); Modernización Sustentable de la Agricultura Tradicional (MasAgro) an initiative from the Secretaría de Agricultura y Desarrollo Rural (SADER), Mexico; Foundation for Food and Agricultural Research (FFAR). **ETHZ**: Swiss National Science Foundation (SNSF) **INRAe**: FFAST (French National Research Agency, ANR project number ANR-21-CE45-0037). **RRes**: Biotechnology and Biological Sciences Research Council (BBSRC) of the UK as part of the project Delivering Sustainable Wheat (BB/X011003/1). **Uliège**: National Fund of Belgium F.R.S-FNRS (FRIA grant), Agriculture, Natural Resources and Environment Research Direction of the Public Service of Wallonia (project D31-1385 PHENWHEAT). **UQ**: INVITA - A technology and analytics platform for improving variety selection, GRDC UOQ2003-011RTX and GRDC UOQ2002-08RTX High-

throughput feature extraction from imagery to map spatial variability. **USask**: the Natural Sciences and Engineering Research Council of Canada (NSERC) and the Canada First Research Excellence Fund (CFREF). **UTokyo**: the Japan Science and Technology Agency AIP Acceleration Research (JPMJCR21U3), the Sarabetsu Village “Endowed Chair for Field Phenomics” project in Hokkaido, Japan.

Data availability

The full dataset (GWFSS_v1.0_full) including the 1096 ground-truth labelled images (GWFSS_v1.0_labelled), the descriptions of the datasets (GWFSS_v1.0_subsets.csv) and imaging setups (GWFSS_v1.0_imaging_setups.csv) is available in the ETH research collection (<https://doi.org/10.3929/ethz-b-000734546>) and should be referred to by citing this publication. The 1096 labelled images will be released on July 15, 2025. To facilitate access, the labelled data and the benchmark model will also be available at (https://huggingface.co/datasets/GlobalWheat/GWFSS_v1.0). Further subsets of the data are available as part of a competition held in the framework of the Seventh International Workshop on Machine Learning for Cyber-Agricultural Systems (MLCAS2025). These sets are listed in Annex 1.1 but are not part of the publication. Links to these datasets can be found at: <https://www.global-wheat.com/gwfss.html>.

Declaration of competing interest

The authors declare that there is no conflict of interest regarding the publication of this article.

Acknowledgments

We thank all field staff who ran the experiments and collected the images used in this study.

Appendix A. Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.plaphe.2025.100084>.

References

- [1] O. Erenstein, M. Jaleta, K.A. Mottaleb, K. Sonder, J. Donovan, H.J. Braun, Global trends in wheat production, consumption and trade, in: *Wheat Improvement: Food Security in a Changing Climate*, Springer International Publishing, Cham, 2022, pp. 47–66.
- [2] N. Brisson, P. Gate, D. Gouache, G. Charmet, F.X. Oury, F. Huard, Why are wheat yields stagnating in Europe? A comprehensive data analysis for France, *Field Crops Res.* 119 (2010) 201–212.
- [3] D.K. Ray, N.D. Mueller, P.C. West, J.A. Foley, Yield trends are insufficient to double global crop production by 2050, *PLoS One* 8 (2013) e66428.
- [4] M. Morisse, D.M. Wells, E.J. Millet, et al., A European perspective on opportunities and demands for field-based crop phenotyping, *Field Crops Res.* 276 (2022) 108371.
- [5] W. Li, S. Wu, W. Wen, et al., Using high-throughput phenotype platform MVS-Pheno to reconstruct the 3D morphological structure of wheat, *AoB Plants* 16 (2024).
- [6] P. Sadeghi-Tehran, N. Virlet, M.J. Hawkesford, A neural network method for classification of sunlit and shaded components of wheat canopies in the field using high-resolution hyperspectral imagery, *Remote Sens.* 13 (2021) 898.
- [7] F.G. Okyere, D. Cudjoe, P. Sadeghi-Tehran, et al., Modeling the spatial-spectral characteristics of plants for nutrient status identification using hyperspectral data and deep learning methods, *Front. Plant Sci.* 14 (2023) 1209500.
- [8] M. Reynolds, S. Chapman, L. Crespo-Herrera, et al., Breeder friendly phenotyping, *Plant Sci.* 295 (2020) 110396.
- [9] S. Madec, K. Irfan, K. Velumani, et al., VegAnn, vegetation annotation of multi-crop RGB images acquired under diverse conditions for segmentation, *Sci. Data* 10 (2023) 302.
- [10] A. Carlier, S. Dandriofosse, B. Dumont, B. Mercatoris, Comparing CNNs and PLSr for estimating wheat organs biophysical variables using proximal sensing, *Front. Plant Sci.* 14 (2023) 1204791.
- [11] N. Virlet, K. Sabermanesh, P. Sadeghi-Tehran, M.J. Hawkesford, Field scanalyzer: an automated robotic field phenotyping platform for detailed crop monitoring, *Funct. Plant Biol.* 44 (2016) 143–153.
- [12] B. Gano, S. Bhadra, J.M. Vilbig, N. Ahmed, V. Sagan, N. Shakoar, Drone-based imaging sensors, techniques, and applications in plant phenotyping for crop breeding: a comprehensive review, *Plant Phenome J.* 7 (2024) e20100.
- [13] K. Yu, N. Kirchgessner, C. Grieder, A. Walter, A. Hund, An image analysis pipeline for automated classification of imaging light conditions and for quantification of wheat canopy cover time series in field phenotyping, *Plant Methods* 13 (2017) 1–13.
- [14] M.M. Hasan, J.P. Chopin, H. Laga, S.J. Miklavcic, Detection and analysis of wheat spikes using convolutional neural networks, *Plant Methods* 14 (2018) 1–13.
- [15] E. David, S. Madec, P. Sadeghi-Tehran, et al., Global wheat head detection (GWHD) dataset: a large and diverse dataset of high-resolution RGB-Labelled images to develop and benchmark wheat head detection methods, *Plant Phenomics* 2020 (2020).
- [16] E. David, M. Serouart, D. Smith, et al., Global wheat head detection 2021: an improved dataset for benchmarking wheat head detection methods, *Plant Phenomics* 2021 (2021).
- [17] S. Madec, X. Jin, H. Lu, et al., Ear density estimation from high resolution RGB imagery using deep learning technique, *Agric. For. Meteorol.* 264 (2019) 225–234.
- [18] R. Zenkl, R. Timofte, N. Kirchgessner, et al., Outdoor plant segmentation with deep learning for high-throughput field phenotyping on a diverse wheat dataset, *Front. Plant Sci.* 12 (2022) 774068.
- [19] M. Serouart, S. Madec, E. David, et al., SegVeg: segmenting RGB images into green and senescent vegetation by combining deep and shallow methods, *Plant Phenomics* 2022 (2022).
- [20] H. Anwar, S.U. Khan, M.M. Ghaffar, et al., The NWRD dataset: an open-source annotated segmentation dataset of diseased wheat crop, *Sensors* 23 (2023) 6942.
- [21] H. Getachew, T. Hawi, Wheat leaf dataset, *Mendeley Data* 1 (2021).
- [22] Y. Li, T. Qiao, W. Leng, et al., Semantic segmentation of wheat stripe rust images using deep learning, *Agronomy* 12 (2022) 2933.
- [23] S. Arya, B. Singh, Wheat nitrogen deficiency and leaf rust image dataset, *Mendeley Data* 1 (2020) 2020.
- [24] S. Hussain, A. Moreno, S. Chauda, CGIAR Computer Vision for Crop Disease, Kaggle, 2021.
- [25] Y. Guan, J. Pan, Q. Fan, L. Yang, X. Yin, W. Jia, CTWheatNet: accurate detection model of wheat ears in field, *Comput. Electron. Agric.* 225 (2024) 109272.
- [26] Y. Wang, F. Wang, K. Li, et al., Low-light wheat image enhancement using an explicit inter-channel sparse transformer, *Comput. Electron. Agric.* 224 (2024) 109169.
- [27] C. Liu, K. Wang, H. Lu, Z. Cao, Dynamic color transform networks for wheat head detection, *Plant Phenomics* 2022 (2022).
- [28] Y. Li, X. Wu, Q. Wang, et al., CSNet: a count-supervised network via multiscale MLP-mixer for wheat ear counting, *Plant Phenom.* 6 (2024) 236.
- [29] S. Dandriofosse, E. Ennadifi, A. Carlier, B. Gosselin, B. Dumont, B. Mercatoris, Deep learning for wheat ear segmentation and ear density measurement: from heading to maturity, *Comput. Electron. Agric.* 199 (2022) 107161.
- [30] Y. Guan, J. Pan, Q. Fan, L. Yang, L. Xu, W. Jia, Generalized focal loss WheatNet (GFLWheatNet): Accurate application of a wheat ear detection model in field yield prediction, *Agriculture* 14 (2024) 899.
- [31] S. Khaki, N. Safaei, H. Pham, L. Wang, WheatNet: a lightweight convolutional neural network for high-throughput image-based wheat head detection and counting, *Neurocomputing* 489 (2022) 78–89.
- [32] K. Najafian, A. Ghanbari, I. Stavness, L. Jin, G.H. Shirdel, F. Maleki, A semi-self-supervised learning approach for wheat head detection using extremely small number of labeled samples, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 1342–1351.
- [33] F. Fourati, W.S. Mseddi, R. Attia, Wheat head detection using deep, semi-supervised and ensemble learning, *Can. J. Rem. Sens.* 47 (2021) 198–208.
- [34] F.C. Ogidi, M.G. Eramian, I. Stavness, Benchmarking self-supervised contrastive learning methods for image-based plant phenotyping, *Plant Phenom.* 5 (2023) 37.
- [35] J. Myers, K. Najafian, F. Maleki, K. Owens, Efficient wheat head segmentation with minimal annotation: a generative approach, *J. Imag.* 10 (2024) 152.
- [36] K. Najafian, A. Ghanbari, M. Sabet Kish, et al., Semi-self-supervised learning for semantic segmentation in images with dense patterns, *Plant Phenom.* 5 (2023) 25.
- [37] C. James, Y. Gu, S. Chapman, et al., Domain adaptation for plant organ detection with style transfer, in: *2021 Digital Image Computing: Techniques and Applications (DICTA)*, IEEE, 2021, pp. 1–9.
- [38] B. Yang, R. Chen, Z. Gao, H. Zhi, FIDMT-GhostNet: a lightweight density estimation model for wheat ear counting, *Front. Plant Sci.* 15 (2024) 1435042.
- [39] P. Sadeghi-Tehran, N. Virlet, E.M. Ampe, P. Reyns, M.J. Hawkesford, DeepCount: in-field automatic quantification of wheat spikes using simple linear iterative clustering and deep convolutional neural networks, *Front. Plant Sci.* 10 (2019) 1176.
- [40] K. Velumani, S. Madec, B. De Solan, et al., An automatic method based on daily in situ images and deep learning to date wheat heading stage, *Field Crops Res.* 252 (2020) 107793.
- [41] P. Sadeghi-Tehran, K. Sabermanesh, N. Virlet, M.J. Hawkesford, Automated method to determine two critical growth stages of wheat: heading and flowering, *Front. Plant Sci.* 8 (2017) 252.
- [42] G.A. Slafer, M.J. Foulkes, M.P. Reynolds, et al., A 'wiring diagram' for sink strength traits impacting wheat yield potential, *J. Exp. Bot.* 74 (2023) 40–71.
- [43] J. Anderegg, N. Kirchgessner, H. Aasen, et al., Thermal imaging can reveal variation in stay-green functionality of wheat canopies under temperate conditions, *Front. Plant Sci.* 15 (2024) 1335037.
- [44] T. Liu, W. Wu, W. Chen, C. Sun, X. Zhu, W. Guo, Automated image-processing for counting seedlings in a wheat field, *Precis. Agric.* 17 (2016) 392–406.
- [45] S. Wang, J. Zhao, Y. Cai, et al., A method for small-sized wheat seedlings detection: from annotation mode to model construction, *Plant Methods* 20 (2024) 15.
- [46] J. Anderegg, F. Tschurr, N. Kirchgessner, et al., Pixel to practice: multi-scale image data for calibrating remote-sensing-based winter wheat monitoring methods, *Sci. Data* 11 (2024) 1033.
- [47] P. Sadeghi-Tehran, N. Virlet, K. Sabermanesh, M.J. Hawkesford, Multi-feature machine learning model for automatic segmentation of green fractional vegetation cover for high-throughput field phenotyping, *Plant Methods* 13 (2017) 1–16.
- [48] S. Brocks, G. Bareth, Estimating barley biomass with crop surface models from oblique RGB imagery, *Remote Sens.* 10 (2018) 268.
- [49] S. Wu, L. Deng, L. Guo, Y. Wu, Wheat leaf area index prediction using data fusion based on high-resolution unmanned aerial vehicle imagery, *Plant Methods* 18 (2022) 68.
- [50] P. Castro-Valdecantos, O.E. Apolo-Apolo, M. Pérez-Ruiz, G. Gege, Leaf area index estimations by deep learning models using RGB images and data fusion in maize, *Precis. Agric.* 23 (2022) 1949–1966.
- [51] C. Grieder, A. Hund, A. Walter, Image based phenotyping during winter: a powerful tool to assess wheat genetic variation in growth response to temperature, *Funct. Plant Biol.* 42 (2015) 387–396.
- [52] G.E. Meyer, J.C. Neto, Verification of color vegetation indices for automated crop imaging applications, *Comput. Electron. Agric.* 63 (2008) 282–293.
- [53] Y. Wang, Z. Cao, X. Bai, Z. Yu, Y. Li, An automatic detection method to the field wheat based on image processing, in: *MIPPR 2013: Automatic Target Recognition and Navigation*, 8918, SPIE, 2013, pp. 97–104.
- [54] L. Roth, M. Camenzind, H. Aasen, et al., Repeated multiview imaging for estimating seedling tiller counts of wheat genotypes using drones, *Plant Phenomics* 2020 (2020).
- [55] F.G. Okyere, D. Cudjoe, P. Sadeghi-Tehran, et al., Machine learning methods for automatic segmentation of images of field-and glasshouse-based plants for high-throughput phenotyping, *PLoS One* 12 (2023) 2035.
- [56] J. Grimm, K. Herzog, F. Rist, A. Kicherer, R. Toepfer, V. Steinhage, An adaptable approach to automated visual detection of plant organs with applications in grapevine breeding, *Biosyst. Eng.* 183 (2019) 170–183.
- [57] S.G. Sodjinou, V. Mohammadi, A.T.S. Mahama, P. Gouton, A deep semantic segmentation-based algorithm to segment crops and weeds in agronomic color images, *Inform. Proc. Agric.* 9 (2022) 355–364.
- [58] A. Milioto, P. Lottes, C. Stachniss, Real-time semantic segmentation of crop and weed for precision agriculture robots leveraging background knowledge in CNNs,

- in: 2018 IEEE International Conference on Robotics and Automation (ICRA), IEEE, 2018, pp. 2229–2235.
- [59] J. Ma, Y. Li, H. Liu, et al., Improving segmentation accuracy for ears of winter wheat at flowering stage by semantic segmentation, *Comput. Electron. Agric.* 176 (2020) 105662.
- [60] K. Zou, Q. Liao, F. Zhang, X. Che, C. Zhang, A segmentation network for smart weed management in wheat fields, *Comput. Electron. Agric.* 202 (2022) 107303.
- [61] R. Zenkl, R. Timofte, N. Kirchgessner, et al., Outdoor plant segmentation with deep learning for high-throughput field phenotyping on a diverse wheat dataset, *Front. Plant Sci.* 12 (2022) 774068.
- [62] H. Peng, C. Xue, Y. Shao, et al., Semantic segmentation of litchi branches using DeepLabV3+ model, *IEEE Access* 8 (2020) 164546–164555.
- [63] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J.M. Álvarez, P. Luo, SegFormer: simple and efficient design for semantic segmentation with transformers, in: *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6–14, 2021, virtual*, 2021, pp. 12077–12090.
- [64] J. Li, M. Xu, L. Xiang, et al., Foundation models in smart agriculture: basics, opportunities, and challenges, *Computers and Electronics in Agriculture* 222 (2024) 109032.
- [65] D. Wang, D. Zhang, G. Yang, B. Xu, Y. Luo, X. Yang, SSRNet: in-Field counting wheat ears using multi-stage convolutional neural network, *IEEE Trans. Geosci. Rem. Sens.* 60 (2021) 1–11.
- [66] J. Ma, Y. Li, K. Du, et al., Segmenting ears of winter wheat at flowering stage using digital images and deep learning, *Comput. Electron. Agric.* 168 (2020) 105159.
- [67] T. Alkhudaydi, D. Reynolds, S. Griffiths, J. Zhou, B. De La Iglesia, An exploration of deep-learning based phenotypic analysis to detect spike regions in field conditions for UK bread wheat, *Plant Phenomics* (2019), 2019.
- [68] T. Alkhudaydi, B. De la Iglesia, Counting spikelets from infield wheat crop images using fully convolutional networks, *Neural Comput. Appl.* 34 (2022) 17539–17560.
- [69] X. Xu, Q. Geng, F. Gao, D. Xiong, H. Qiao, X. Ma, Segmentation and counting of wheat spike grains based on deep learning and textural feature, *Plant Methods* 19 (2023) 77.
- [70] A. Carlier, S. Dandriofosse, B. Dumont, B. Mercatoris, To what extent does yellow rust infestation affect remotely sensed nitrogen status? *Plant Phenom.* 5 (2023) 83.
- [71] T. Liu, Y. Zhao, Y. Sun, et al., High-throughput identification of fusarium head blight resistance in wheat varieties using field robot-assisted imaging and deep learning techniques, *J. Clean. Prod.* 480 (2024) 144024.
- [72] J. Deng, X. Lv, L. Yang, et al., Assessing macro disease index of wheat stripe rust based on segformer with complex background in the field, *Sensors* 22 (2022) 5676.
- [73] J. Anderegg, R. Zenkl, A. Walter, A. Hund, B.A. McDonald, Combining high-resolution imaging, deep learning, and dynamic modeling to separate disease and senescence in wheat canopies, *Plant Phenom.* 5 (2023) 53.
- [74] B de Solan, G. Daubige, S. Thomas, M. Serouart, LITERAL, a portable, lightweight and accurate high-throughput phenotyping system for crop monitoring, *Innovat. Agronom.* 94 (2024) 231–240.
- [75] N. Kirchgessner, F. Liebisch, K. Yu, et al., The ETH field phenotyping platform FIP: a cable-suspended multi-sensor system, *Funct. Plant Biol.* 44 (2016) 154–168.
- [76] L. Roth, M. Boss, N. Kirchgessner, et al., The FIP 1.0 data set: highly resolved annotated image time series of 4,000 wheat plots grown in 6 years, *GigaScience* (2025), 14:giaf051.
- [77] W. Guo, T. Fukatsu, S. Ninomiya, Automated characterization of flowering dynamics in rice using field-acquired time-series RGB images, *Plant Methods* 11 (2015) 7.
- [78] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [79] A. Richardson, H. Jones, M. Bartlett, Grass awns: morphological diversity arising from developmental constraint, *Curr. Opin. Plant Biol.* 82 (2024) 102663.
- [80] L. Chen, Y. Zhu, G. Papandreou, F. Schroff, H. Adam, Encoder-decoder with atrous separable convolution for semantic image segmentation, in: V. Ferrari, M. Hebert, C. Sminchisescu, Y. Weiss (Eds.), *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8–14, 2018, Proceedings, Part VII*, Lecture Notes in Computer Science, 11211, Springer, 2018, pp. 833–851, https://doi.org/10.1007/978-3-030-01234-2%5C_49, 10.1007/978-3-030-01234-2_49.
- [81] L. Roth, A. Hund, H. Aasen, PhenoFly planning tool: flight planning for high-resolution optical remote sensing with unmanned areal systems, *Plant Methods* 14 (2018) 1–21.
- [82] E.A. Papoutsoglou, D. Faria, D. Arend, et al., Enabling reusability of plant phenomic datasets with MIAPPE 1.1, *New Phytol.* 227 (2020) 260–273.
- [83] C. Celestina, J. Hunt, H. Brown, et al., Scales of development for wheat and barley specific to either single culms or a population of culms, *Eur. J. Agron.* 147 (2023) 126824.
- [84] J.C. Zadoks, T.T. Chang, C.F. Konzak, A decimal code for the growth stages of cereals, *Weed Res.* 14 (1974) 415–421.
- [85] P.D. Lancashire, H. Bleiholder, Tvd Boom, et al., A uniform decimal code for growth stages of crops and weeds, *Ann. Appl. Biol.* 119 (1991) 561–601.
- [86] R. Kestur, A. Meduri, O. Narasipura, MangoNet: a deep semantic segmentation architecture for a method to detect and count mangoes in an open orchard, *Eng. Appl. Artif. Intell.* 77 (2019) 59–69.
- [87] N. Chebrolu, P. Lottes, A. Schaefer, W. Winterhalter, W. Burgard, C. Stachniss, Agricultural robot dataset for plant classification, localization and mapping on sugar beet fields, *Int. J. Robot Res.* 36 (2017) 1045–1052.
- [88] Y. Lu, S. Young, A survey of public datasets for computer vision tasks in precision agriculture, *Comput. Electron. Agric.* 178 (2020) 105760.
- [89] Z. Niu, N. Liang, Y. He, et al., A novel method for wheat spike phenotyping based on instance segmentation and classification, *Appl. Sci.* (2024) 14.
- [90] R. Zenkl, B.A. McDonald, A. Walter, M. Unvericht, C. Saintenac, J. Anderegg, From Canopy Images to Organ-Level Disease Assessments: A Scalable Approach to Measure Quantitative Resistance in the Field, *bioRxiv* (2025), 2025–4.
- [91] E.A. Chapman, S. Orford, J. Lage, S. Griffiths, Capturing and selecting senescence variation in wheat, *Front. Plant Sci.* 12 (2021) 638738.
- [92] C. James, D. Smith, W. He, S.S. Chandra, S.C. Chapman, GrainPointNet: a deep-learning framework for non-invasive sorghum panicle grain count phenotyping, *Comput. Electron. Agric.* 217 (2024) 108485.
- [93] D.M. Deery, G.J. Rebetzke, J.A. Jimenez-Berni, et al., Ground-based LiDAR improves phenotypic repeatability of above-ground biomass and crop growth rate in wheat, *Plant Phenomics*, 2020.
- [94] S. El-Hendawy, M.U. Tahir, N. Al-Suhaibani, S. Elsayed, O. Elsherbiny, H. Elsharawy, Potential of thermal and RGB imaging combined with artificial neural networks for assessing salt tolerance of wheat genotypes grown in real-field conditions, *Agronomy* 14 (2024) 1390.
- [95] Tools and Methods for Extended Plant Phenotyping and EnviroTyping Services of European Research Infrastructures, 2025.
- [96] Bringing EMPHASIS to Operation: European Infrastructure for multi-scale Plant Phenomics and Simulation for Food Security in a Changing Climate, 2022.