

Transformer-based Detection of Abnormal Rice Growth using Drone-based Multispectral Imaging

Yanfen Li^{a,1}, L. Minh Dang^{b,c,d,1}, Hanxiang Wang^{a,*}, Muhammad Fayaz^e,
Sufyan Danish^e, Junliang Shang^a, Hyoung-Kyu Song^d, Hyeonjoon Moon^{e,*}

^a*School of Computer Science, Qufu Normal University, Rizhao, China*

^b*Institute of Research and Development, Duy Tan University, Da Nang, 550000, Viet Nam*

^c*Faculty of Information Technology, Duy Tan University, Da Nang, 550000, Viet Nam*

^d*Department of Information and Communication Engineering and Convergence*

Engineering for Intelligent Drone, Sejong University, Seoul 05006, South Korea

^e*Department of Computer Science and Engineering, Sejong University, Seoul 05006, South Korea*

Abstract

Rice is a vital staple food for global food security and a primary income source for millions of farmers worldwide. However, abnormal rice growth poses a serious threat to both yield stability and grain quality, undermining agricultural productivity. Early detection of such anomalies is therefore essential to mitigate yield losses. However, existing methods either targeted only one symptom at a time, or failed to generalize under various field conditions. Moreover, lightweight real-time inference is needed for on-board UAV deployment, yet most high-accuracy models incur prohibitive computational cost. In this study, we propose ARG-TR model, a lightweight transformer-based semantic segmentation framework built on the SegFormer architecture, which utilizes long-range dependencies to identify complex growth anomalies. The model is trained and validated on a large-scale, drone-captured multi-spectral dataset. By integrating a hierarchical transformer encoder with a lightweight decoder, ARG-TR achieves rapid convergence during training and demonstrates strong generalization to unseen data. The experimental results on a challenging dataset of abnormal rice growth

*Corresponding authors

Email addresses: hanxiang@qufu.edu.cn (Hanxiang Wang), hmoon@sejong.ac.kr (Hyeonjoon Moon)

¹These authors contributed equally to this work.

patterns show that ARG-TR achieves a robust Intersection over Union (IoU) of 64.8, which outperforms state-of-the-art baselines such as MaskFormer and KNet in both accuracy and computational efficiency.

Keywords: abnormal growth, rice, transformers, semantic segmentation, lodging

1. Introduction

The Food and Agriculture Organization (FAO) of the United Nations (Food & of the United Nations, 2025) projects that the global population will reach 9.2 billion by 2050. To meet the food demands of this growing population, global agricultural production must increase by 60–70% from current levels, as emphasized in multiple FAO reports (Samal et al., 2022; Stankus, 2021). Rice, a staple food for over half the world’s population, predominantly in Asia, plays a critical role in global food security (Bin Rahman & Zhang, 2023). However, it is increasingly challenging to achieve the required production target due to abnormal growth patterns in rice, which manifest through various symptoms including stunted growth, delayed flowering, malformed grains, lodging, missing plants, and disease-specific damages, such as rice blast disease. These abnormalities stem from biotic (e.g., diseases, pests) and abiotic stressors (e.g., nutrient deficiencies, environmental pressures) (Dang et al., 2024), which disrupt normal crop development by reducing photosynthetic efficiency and compromising plant components (Rezvi et al., 2023). For example, rice blast disease (*Magnaporthe oryzae*) destroys photosynthetic tissues, while pest infestations weaken vital components, directly compromising both yield quantity and quality. These issues threaten to destabilize rice production if left untreated, undermining food security, farmers’ livelihoods, and economic stability in rice-dependent regions. Therefore, it is urgent to address abnormal rice growth through targeted mitigation strategies to safeguard sustainable rice production and global food security.

Traditionally, abnormal rice growth detection and diagnosis relies heavily on manual field inspections by agricultural experts. However, this approach is

labor-intensive, time-consuming, and impractical for large-scale monitoring due to the massive size of rice fields. Inspecting individual rice plants for subtle growth anomalies on vast areas is logistically unfeasible. To address these challenges, automated inspection systems are critical for enabling timely, field-scale assessment of crop health. Unmanned aerial vehicles (UAVs), equipped with RGB, multispectral, or thermal imaging sensors provide a viable solution for early-stage anomaly detection (Dang et al., 2020). By capturing high-resolution aerial imagery, UAVs offer a bird’s-eye view that reveals subtle stress indicators, such as chlorosis, stunted growth, or canopy structural variations, often undetectable at ground level. UAVs generate large-scale datasets that motivate the development of advanced computer vision (CV) frameworks capable of automated feature extraction, anomaly classification, and quantifiable stress mapping to transform raw imagery into actionable interventions.

While conventional ML methods depend on handcrafted feature extraction, deep learning (DL) models, particularly convolutional neural networks (CNNs), demonstrate superior capacity for automated detection of subtle rice growth abnormalities from UAV or ground-level imagery (Alam et al., 2025). DL models automatically learn discriminative features, such as color, texture, structural, and spatial domains, to identify issues such as stunted growth, disease, or nutrient deficiencies with minimal human intervention. As a result, DL has achieved state-of-the-art performance in various tasks for precision agriculture, including classification (Li et al., 2020; Dang et al., 2020), detection (Dosset et al., 2025; Wang et al., 2024), and segmentation (Alam et al., 2025; Zhang et al., 2021a). For example, Tian et al. (Tian et al., 2021) employed partial least squares discrimination analysis on UAV multispectral data to detect rice lodging. By utilizing spectral, textural, and color features, the model achieved over 90% accuracy. However, its handcrafted spectral features exhibited limited generalization for various cultivars, growth stages, and regional conditions because the model was fine-tuned on Shanghai paddy characteristics. With a more advanced approach, Yang et al. (Yang et al., 2020) introduced an adaptive UAV-based scouting system that combines multi-altitude imaging and a deep segmentation

56 model to detect rice and lodging. The model achieved 95.28% rice identifica-
 57 tion and 86.17% lodging detection. However, the results are based primarily
 58 on simulations and selected UAV energy profiles. On the other hand, Zhang
 59 et al. (Zhang et al., 2021a) developed Ir-UNet, a DL model for wheat yellow
 60 rust detection. By integrating irregular convolution and content-aware channel
 61 reweighting modules, Ir-UNet addressed challenges posed by irregularly shaped
 62 and blurred disease boundaries. The experimental results showed that the model
 63 achieved 97.13% overall accuracy on UAV multispectral data and maintained
 64 robustness with reduced input features. Recently, Wu et al. (Wu et al., 2025)
 65 proposed a YOLOv5-based pipeline for missing rice seedling detection using
 66 UAV images. UAV images were first stitched into a geo-referenced panoramic
 67 view and then cropped to a series 640×640 patches for dataset creation. The
 68 patches were used to train a YOLOv5, which achieved an 80% recall and 75%
 69 precision in identifying missing rice seedlings. However, GPS-dependent image
 70 stitching and predefined thresholds degraded performance in fields with irregu-
 71 lar planting patterns or GPS drift, and the rectangular detection regions could
 72 miss seedlings in non-uniform layouts. In general, previous approaches suffer
 73 from three main limitations: (1) single-symptom detection, such as lodging or
 74 single disease detection, (2) poor generalizability due to limited labeled training
 75 data, (3) limited adaptability to irregular input due to grid layouts.

76 Originally developed for natural language processing (NLP), transformer
 77 models revolutionized CV by introducing a self-attention mechanism to model
 78 long-range spatial dependencies and global context (Lin et al., 2022). The Vision
 79 Transformer (ViT) (Dosovitskiy et al., 2020) pioneered this for CV by partition-
 80 ing images into patch tokens, but its computational inefficiency limited dense
 81 prediction tasks. Swin Transformer (Liu et al., 2021) addressed this by intro-
 82 ducing hierarchical feature extraction and shifted windowing scheme to improve
 83 efficiency and spatial reasoning for dense prediction tasks like semantic segmen-
 84 tation. Building on these innovations, SegFormer (Xie et al., 2021) emerged as
 85 a state-of-the-art semantic segmentation model. It combined transformer-based
 86 global context modeling with a lightweight, hierarchical architecture to simul-

87 taneously capture fine-grained details and broader contextual relationships. It
 88 achieved high accuracy of 84.0% on Cityscapes with only around 3.8 million
 89 parameters. Therefore, SegFormer was proved to be suitable for tasks requiring
 90 precise localization of subtle anomalies like abnormal rice growth identification.

91 Building on SegFormer’s efficiency and robustness in agricultural applica-
 92 tions (Spasev et al., 2024; Nuradili et al., 2024), this study proposes a light-
 93 weight transformer-based framework engineered to overcome the multi-symptom
 94 detection gap identified in Section 2. The model simultaneously identifies four
 95 different rice growth abnormalities using multispectral imaging. Key contribu-
 96 tions include.

- 97 • Comprehensive data processing and effective post-processing to generate
 98 precise orthophotos for the collected multi-spectral dataset.
- 99 • A large-scale UAV-based remote sensing dataset containing over 378,000
 100 images.
- 101 • An optimized spectral fusion of green, near-infrared, and red-edge to im-
 102 prove image quality for accurate abnormal rice growth recognition.
- 103 • A light-weight transformer-based system for the identification of four ab-
 104 normal rice growth symptoms.

105 The rest of this paper is organized as follows. Section 3 describes the ab-
 106 normal rice growth dataset used in this study. Section 4 presents the proposed
 107 ARG-TR framework for multi-spectral rice growth segmentation. Section 5 re-
 108 ports the experimental setup and results. Section 6 discusses the main findings,
 109 limitations, and practical implications. Finally, Section 7 concludes the paper
 110 and outlines directions for future research.

111 2. System Overview

112 Figure 1 depicts the main processes of AGR-TR, a multi-symptom abnormal
 113 rice growth segmentation framework. In this context, AGR indicates that the

framework is applied to agriculture context, while TR refers to the transformer-based architecture.

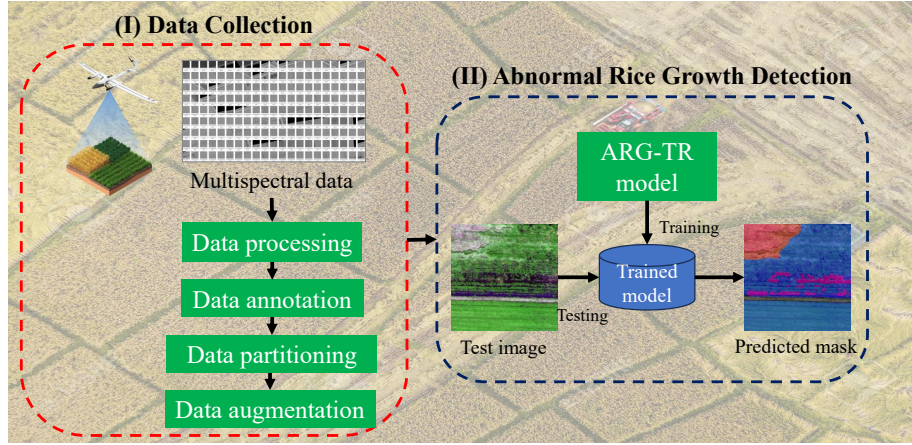


Figure 1: Depiction of the main components of the transformer-based abnormal rice growth detection framework (AGR-TR).

The framework consists of two sequential stages: data preparation and abnormal rice growth detection. In the first stage, multi-spectral UAV imagery undergoes various preprocessing steps to reduce noise and enhance quality. The preprocessed images are then annotated with pixel-level labels to distinguish abnormal growth regions. Next, the dataset is partitioned into training and validation sets. Finally, the images from the training set is augmented to improve model robustness. The second stage employs and fine-tunes a transformer-based SegFormer to segment abnormal growth areas. The model is trained on the processed dataset to automatically discriminate against spatial-spectral patterns. During inference, the trained model processes an input image to generate an output mask that highlights regions of abnormal rice growth. This end-to-end pipeline integrates advanced CV techniques with agronomic insights to support real-time rice health monitoring.

129 3. Abnormal Rice Growth Dataset

130 This study utilizes a large-scale abnormal rice growth dataset comprising
131 approximately 378,000 multi-spectral images capturing four distinct patterns of
132 abnormal growth. The dataset was made available for research purposes by the
133 National Information Society Agency of Korea (NIA)², which ensures robustness
134 and practical relevance for real-world agricultural applications. The dataset
135 was developed through a collaborative initiative led by Geomatic Limited³ in
136 partnership with various organizations. Sunyoungeng Limited⁴ manages data
137 collection, whereas NEWLAYER Limited⁵ and Muhanit Limited⁶ handles data
138 annotation and preprocessing.

139 3.1. Data collection

140 Abnormal rice growth data were collected from 2022 to 2023 in a 100-hectare
141 experimental crop field in Jangan-ri, Jangan-myeon, Hwaseong City, Gyeonggi
142 Province, South Korea (Figure 2). The site features a temperate monsoon cli-
143 mate ideal for rice cultivation, characterized by warm, humid summers (25–30
144 °C) and annual rainfall of 1,100–1,400 mm concentrated during the summer
145 monsoon season. Fertile loamy to clay-loamy soils (pH 5.5–7.0) ensure strong
146 water retention and nutrient availability (Ju et al., 2022), while the flat topog-
147 raphy supports efficient irrigation and uniform field management.

148 The *Oryza sativa* 'Odae' cultivar (a widely cultivated Japonica variety) was
149 transplanted on 26 May 2022 at a density of 30 cm × 17 cm. Fertilization fol-
150 lowed regional standards with applications of nitrogen (89 kg/hm²), phospho-
151 rus (40 kg/hm²), and potash (53 kg/hm²). Data collection covered five critical
152 growth stages, including tillering, panicle initiation, booting, heading & flow-
153 ering, grain filling. This study specifically targets four high-impact abnormal

²https://www.nia.or.kr/site/nia_kor/main.do

³<https://www.geomatic.co.kr/>

⁴http://nonghyup.ac.kr/e_main.asp

⁵<http://egis.everlinks.co.kr/>

⁶<https://muhanit.kr/>

154 growth groups: (1) missing plants (indicating seedling establishment failure), (2)
 155 lodging (stem collapse compromising harvest efficiency), (3) rice blast disease
 156 (*Magnaporthe oryzae* infection causing necrotic lesions), and (4) poor growth
 157 (exhibiting chlorosis, reduced tillering, and diminished vigor).

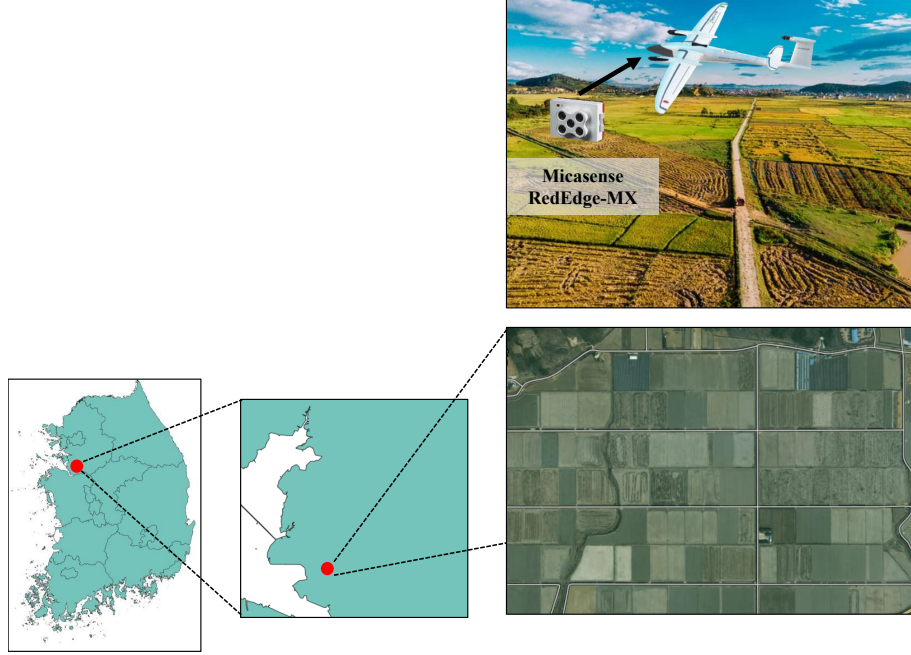


Figure 2: Abnormal rice growth test bed.

158 The 100 hectares experimental rice field was divided into 40 zones, with data
 159 collected from 10 representative plots per zone, leading to a total of 400 moni-
 160 tored plots. Data acquisition main focus was on capturing high-resolution RGB
 161 and multispectral imagery to identify rice field abnormalities using TRINITY
 162 F90+ (Measurusa, 2025). The TRINITY F90+ is a certified vertical take-off
 163 and landing mapping drone, which features a 2.394 m wingspan and a 5.0 kg
 164 maximum take-off weight. With a maximum flight time of 90 minutes and op-
 165 erational range of up to 100 km, it can cover approximately 700 hectares in a
 166 single flight. The drone is compatible with various payloads, including the Mi-
 167 casense RedEdge-MX multispectral camera (MicaSense, 2025), which captures

168 imagery on five narrow spectral bands (blue, green, red, red-edge, and near-
169 infrared). The integration of RedEdge-MX with the TRINITY F90+ enabled
170 efficient, high-fidelity multispectral data acquisition and provided detailed in-
171 sights into crop health, stress detection, and growth dynamics. Flights were
172 performed at 120 m altitude and 5 m/s ground speed, which achieved a ground
173 sample distance (GSD) of 8 cm per pixel.

174 Environmental variables, such as wind gusts, lighting conditions and phe-
175 nological factors, can significantly influence spectral interpretation. For exam-
176 ple, midday sun creates strong shadows that exaggerate canopy gaps, while
177 variable solar angles alter reflectance baselines for chlorosis detection. On the
178 other hand, late-season tillering changes the reflectance baseline against which
179 stunting or chlorosis is detected. To address environmental variables affecting
180 spectral interpretation, the data acquisition implemented three critical controls:
181 (1) All flights conducted between 09:00-11:00 KST under clear-sky conditions
182 to minimize solar angle variation, (2) Geometric correction using calibration
183 and ground control point, (3) collection of five growth stages (tillering to grain
184 filling) to enable robustness against canopy architecture changes.

185 *3.2. Data processing*

186 Figure 3 illustrates the end-to-end workflow for analyzing abnormal rice
187 growth using drone-based multispectral imagery from the experimental field.
188 Initially, all raw imagery undergoes internal quality assurance review by the
189 lead data collector before transmission to the processing team. The workflow
190 begins with raw multispectral data, which is aligned to ensure consistent spa-
191 tial overlap between images. Next, spectral calibration corrects environmental
192 variability (e.g., lighting, atmospheric conditions) and sensor inconsistencies.
193 Ground control point (GCP) correction then enhances positional accuracy by
194 aligning image coordinates with real-world locations (Agüera-Vega et al., 2017),
195 followed by geometric correction to address distortions from sensor tilt or terrain
196 variations. These preprocessing steps collectively produce precise, georeferenced
197 orthophotos, used for subsequent analysis.

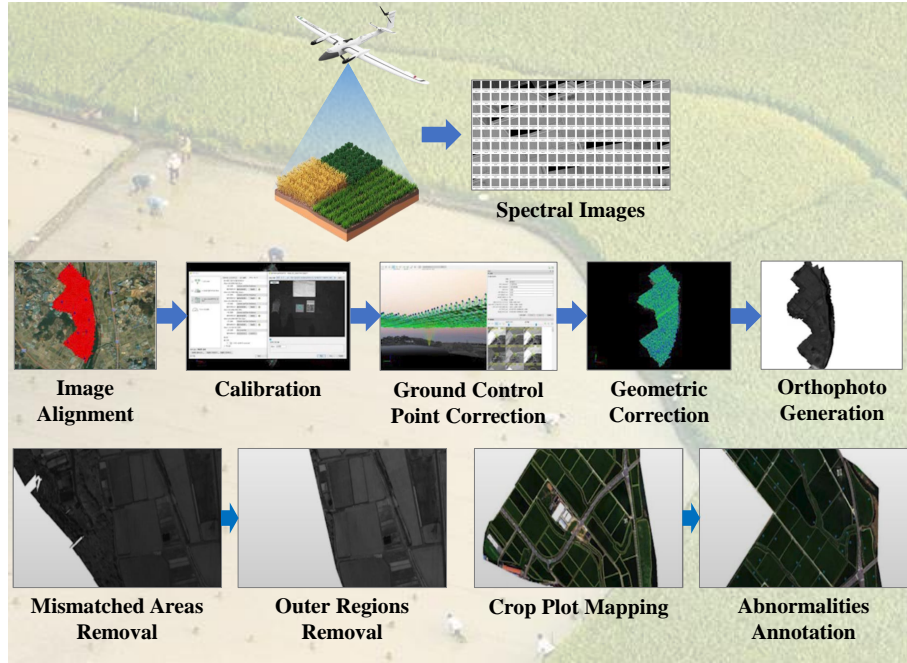


Figure 3: Depiction of the main processing steps for the collected abnormal rice growth dataset.

198 Prior research has established the critical role of Green, Near-infrared (NIR),
 199 and Red-edge (RE) spectral bands for vegetation analysis (Biswal et al., 2024;
 200 Kang et al., 2021). Biswal et al. (Biswal et al., 2024) demonstrated the exclu-
 201 sive use of these three bands for estimating paddy crop aboveground biomass,
 202 while Kang et al. (Kang et al., 2021) highlighted the role of features derived
 203 from RE-NIR-Green band combinations in crop classification. Building on this
 204 foundation, Green, NIR, and RE bands were merged to segment abnormal rice
 205 growth, as the merged version enhanced detection of plant stress, growth anoma-
 206 lies, and terrain characteristics (Dang et al., 2024). Figure 4 illustrates the
 207 creation of combined RGB-like images from these bands within multispectral
 208 orthophotos. This process merges individual channel images into a 3-channel
 209 format compatible with standard CV algorithms and DL models.

210 Finally, a post-processing pipeline was carried out to ensure the usability and

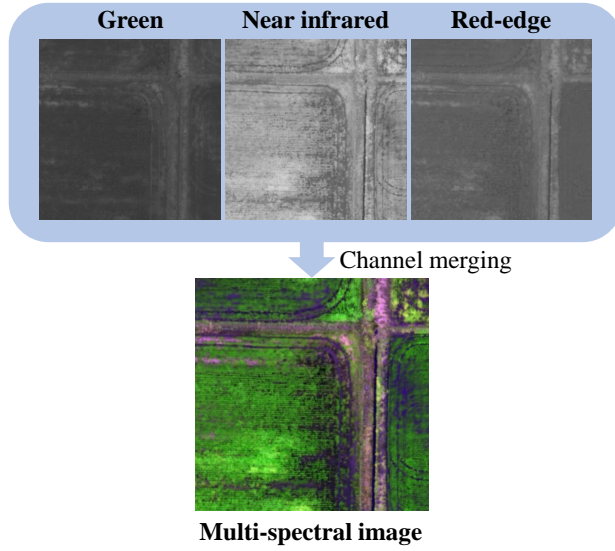


Figure 4: Example of RGB-like image generation through channel merging .

211 accuracy of the dataset. Irrelevant or distorted sections, such as mismatched
 212 areas and outer regions, were removed. The refined orthophoto was divided
 213 into crop plots for localized analysis. Finally, abnormalities, such as rice blast
 214 disease, lodging, poor growth, and missing plants, were labeled to support model
 215 training and validation.

216 The drone’s high-resolution camera delivered a GSD of 8 cm/pixel, sufficient
 217 to resolve individual plants and small abnormal growth patches. For annotation
 218 process, we overlaid the farm-plot map onto the orthomosaics and annotated
 219 each plot showing abnormal growth. To ensure the accuracy of annotations
 220 in drone imagery captured at 120 meters altitude, a multi-stage verification
 221 protocol was implemented. Vegetation indices, including Normalized Difference
 222 Vegetation Index (NDVI)/Enhanced Vegetation Index (EVI), were computed to
 223 highlight potential anomalies invisible in visible spectra. For example, missing
 224 plants were identified by marking regions with NDVI values below a predefined
 225 threshold, while poor growth was annotated using EVI. To reduce ambiguity
 226 and improve annotation consistency, the missing plants class is strictly defined

227 as continuous bare-soil regions within the planted area rather than isolated gaps.
 228 In practice, annotators marked a region as missing plants only when the bare-
 229 soil patch formed a coherent area spanning multiple adjacent planting rows or
 230 otherwise appeared as a continuous discontinuity in the crop canopy. Single-
 231 plant gaps, isolated pixels, thin shadows, wheel tracks, and other narrow non-
 232 crop features were annotated as healthy crop. Lodging was labeled by converting
 233 imagery to RGB format and marking the flattened rice locations. However, rice
 234 blast disease showed no distinctive spectral signatures that could be initially
 235 labeled. Therefore, it was labeled immediately post-flight by trained crowd
 236 workers using handheld RGB cameras and GPS markers during field surveys.
 237 All annotators completed rigorous training in multispectral image interpretation
 238 prior to labeling. Moreover, annotated images were validated through random
 239 field surveys to confirm the presence of annotated abnormalities. This integrated
 240 approach ensured annotation reliability despite the challenges of high-altitude
 241 aerial observation.

242 The drone-based multispectral imaging approach offers valuable insights into
 243 abnormal rice growth but faces several data acquisition limitations. (1) Oper-
 244 ational constraints such as limited drone flight time, altitude restrictions, and
 245 narrow camera field of view complicated data collection and processing. (2)
 246 Variations in solar illumination required complicated post-collection data pro-
 247 cessing, and data were collected only at five discrete growth stages with 7-10
 248 day intervals, potentially missing rapid rice blast disease developments. (3)
 249 The study’s focus on a single geographic location with specific soil and climate
 250 conditions limits generalizability to other regions.

251 *3.3. Dataset description*

252 Figure 5 presents the class distribution of the abnormal rice growth dataset.
 253 The dataset contains 378,074 annotated images in five classes: normal condi-
 254 tions, rice blast disease, lodging, poor growth, and missing plants. For model
 255 development, 80% of the dataset (302,450 images) was used for training and val-
 256 idation purposes (226,853 images for training and 75,597 images for validation),

257 20% of the dataset was used for testing (75,624 images).

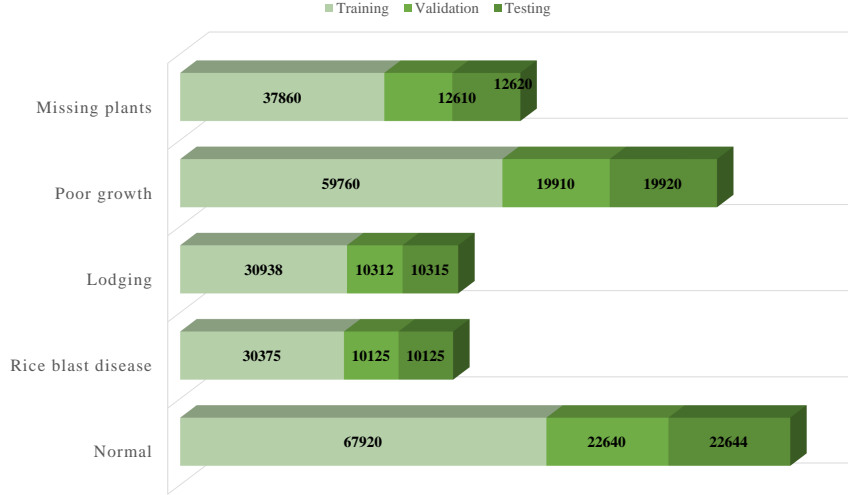


Figure 5: A bar chart showing the distribution of images across different classes in the collected dataset, including normal conditions, rice blast disease, lodging, poor growth, and missing plants.

258 3.4. Data augmentation

259 Data augmentation plays a vital role in enhancing model robustness for
 260 rare anomalies, including rice blast disease, lodging, and missing plants, by
 261 mitigating class imbalance in the training set. A multi-stage augmentation
 262 pipeline was implemented to improve generalization in spatial, spectral, and
 263 scale variation. Figure 6 provides examples of augmented images using the
 264 augmentation pipeline.

265 The images were first scaled by a random factor ranging from 0.5 to 2.0
 266 followed by resizing to 512×512 pixels. This step aims to improve the model
 267 multi-scale robustness by simulating variations in object scale and distance. Af-
 268 ter that, a RandomCrop operation was implemented to sample various regions
 269 to increase diversity in spatial composition. Next, the images were flipped ran-
 270 domly to introduce invariance to orientation. Finally, color jittering (brightness,
 271 contrast, saturation) was applied to mimic diverse lighting conditions and sen-

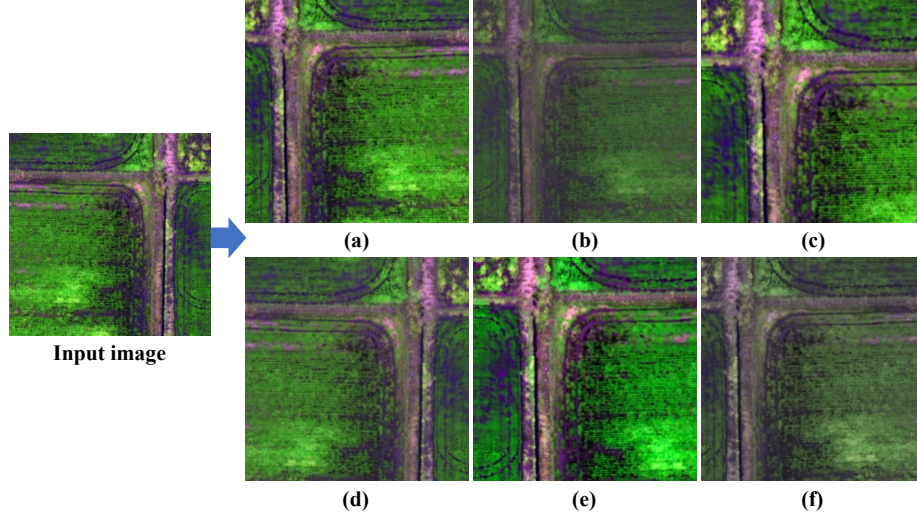


Figure 6: Examples of augmented images (a-f) used for training the ARG-TR model. Augmentations include a combination of scaling, cropping, flipping, and photometric adjustments.

272 sor variations. The augmentation techniques expanded the original training set
273 of 226,853 images by three-fold to 680,550 images.

274 The large-scale and diverse nature of the dataset in capturing multiple types
275 of rice growth abnormalities in different growth stages and environmental con-
276 ditions presented both opportunities and challenges for analysis. The need to
277 effectively process high-resolution multispectral imagery while accurately seg-
278 menting and classifying various abnormal growth patterns in real-time prompted
279 the authors to choose a light-weight framework, which utilizes the rich spectral
280 information in the dataset through self-attention mechanisms while maintaining
281 computational efficiency without sacrificing scalability.

282 4. Methodology

283 Figure 3 illustrates the ARG-TR framework, a Transformer-based system for
284 detecting and segmenting abnormal rice growth. In this context, ARG denote

Abnormal Rice Growth, while TR refers to a light-weight Transformer-based segmentation model.

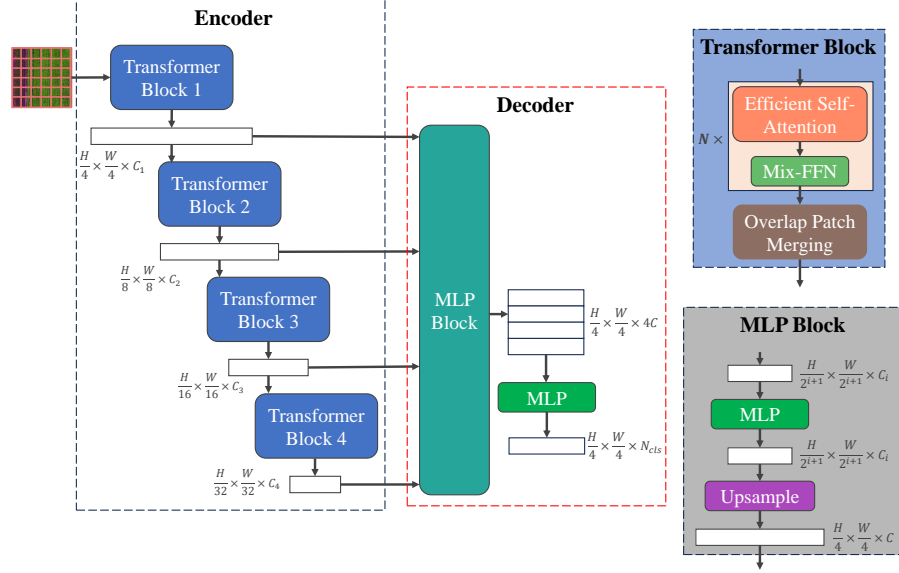


Figure 7: Schematic overview of the ARG-TR framework for abnormal rice growth segmentation. Figure adapted from (Xie et al., 2021)

SegFormer (Xie et al., 2021) is an efficient semantic segmentation architecture that combines a hierarchical transformer encoder and a lightweight multi-layer perceptron (MLP) decoder. Unlike CNN-based models, SegFormer eliminates positional embeddings through overlapped patch merging, which enables consistent performance on variable input sizes while preserving computational efficiency. This study utilizes SegFormer Mix Transformer(MiT)-b3 variant as the foundation. Its key innovations include:

- **Multi-scale encoder:** The encoder extracts both coarse and fine-grained features at four resolutions ($1/4$, $1/8$, $1/16$, and $1/32$ input scale) via overlapping 4×4 patches. Unlike traditional approaches, it does not require positional embeddings to progressively capture fine details and contextual semantics.

- Efficient decoder: The decoder aggregates multi-scale features through MLP layers and upsamples them to produce a high-resolution segmentation map. The decoder ensures precise localization of abnormal growth patterns by fusing coarse (contextual) and fine-grained (detail-rich) features. Channel dimensionality is reduced from 1,024 to 128 via MLP blocks before generating 5-class logits (normal, blast, lodging, poor growth, missing plants).

Table 1 describe the detailed network structure of the ARG-TR model. It begins with a 7×7 convolutional patch embedding (stride=4) to downsample the input to $H/4 \times W/4$ resolution with 64 channels, followed by LayerNorm for normalization and GELU for non-linearity. The encoder consists of four hierarchical stages of Transformer layers: Stage 1 with 3 layers (64 channels, $H/4 \times W/4$), Stage 2 with 3 layers (128 channels, $H/8 \times W/8$), Stage 3 with 18 layers (320 channels, $H/16 \times W/16$), and Stage 4 with 3 layers (512 channels, $H/32 \times W/32$). In the decoder, features from all encoder stages are upsampled to $H/4 \times W/4$, concatenated, and processed through an MLPBlock that reduces the channel dimension from 1024 to 256 to 128, followed by a 1×1 convolution to generate 5 class logits. The head applies a softmax operation to convert logits into probabilities and resizes the output to the original image resolution. This design efficiently captures both fine and coarse details for multispectral rice growth anomaly detection.

4.1. Hierarchical Transformer Encoder

The Mix Transformer (MiT) backbone in SegFormer (Xie et al., 2021) serves as a hierarchical encoder customized for efficient semantic segmentation. It implements a four-stage pyramid structure to generate multi-scale feature maps at resolutions of $1/4$, $1/8$, $1/16$, and $1/32$ of the input image. This design enables robust segmentation of objects for varying scales, from fine-grained details to broader contextual patterns. Between transformer layers, a Mixed Feed-Forward Network (Mix-FFN) integrates depthwise 3×3 convolutions with standard MLP

Table 1: Network structure of the ARG-TR

Module	Layer / Operation	Channels	Output Size
Patch Embedding	Conv 7×7 , stride 4	64	$H/4 \times W/4$
	LayerNorm + GELU	—	$H/4 \times W/4$
Encoder	3 layers (Stage 1)	64	$H/4 \times W/4$
	3 layers (Stage 2)	128	$H/8 \times W/8$
	18 layers (Stage 3)	320	$H/16 \times W/16$
	3 layers (Stage 4)	512	$H/32 \times W/32$
Decoder	Upsampling	—	$H/4 \times W/4$
	Concat \rightarrow MLPBlock	$1024 \rightarrow 256 \rightarrow 128$	$H/4 \times W/4$
	1×1 Conv \rightarrow 4 logits	$128 \rightarrow 4$	$H/4 \times W/4$
Head & Loss	Softmax + resize	—	$H/4 \times W/4 \times 5$ classes
	IoU	—	—

328 operations to enhance local spatial feature interactions. Finally, a patch merging
 329 module downsamples feature maps by concatenating neighboring patches and
 330 linearly projecting channel dimensions to establish a coarse-to-fine feature hi-
 331 erarchy. Moreover, overlapped patch embedding in early stages maintains local
 332 continuity without the need for positional encodings.

333 4.1.1. Hierarchical Feature Representation

334 SegFormer’s encoder generates multi-scale feature maps at $(1/4, 1/8, 1/16,$
 335 $1/32$ input spatial resolution), a crucial improvement from traditional ViTs,
 336 which produce single-scale representations. This hierarchical structure enables
 337 high-resolution feature maps to capture fine-grained details (early-stage rice
 338 blast lesions), while low-resolution feature maps encode coarse contextual in-
 339 formation (lodging propagation). For rice growth analysis, hierarchical feature
 340 representation is essential as fine-grained features detect subtle spectral devia-
 341 tions in individual plants, whereas coarse features model spatial dependencies
 342 across field conditions.

343 4.1.2. Overlapped Patch Merging (OPM)

344 Overlapped Patch Merging (OPM) is an important component of SegFormer’s
 345 Mix Transformer (MiT) encoder that enables hierarchical feature extraction
 346 while preserves local spatial continuity. Unlike standard ViTs with non-overlapping
 347 patches, OPM generates overlapping patches to maintain fine-grained spatial re-
 348 lationships essential for segmenting subtle abnormalities.

349 Given multi-spectral input $\mathbf{X} \in \mathbb{R}^{H \times W \times C}$, where $C = 3$ for Green/NIR/RE
 350 bands, H and W are the height and width. OPM slides a patch window across
 351 the image with a stride S smaller than the patch size K and with padding P ,
 352 so adjacent patches overlap. The stride S being smaller than the patch size K
 353 is the key element that creates the overlap and shared context. Each window
 354 is flattened and linearly projected to form a token for the next hierarchical
 355 level. Repeating this merging yields hierarchical feature maps whose spatial
 356 resolution is reduced (for example from $H/4 \times W/4$ to $H/8 \times W/8$) while the
 357 channel dimension increases. The overlap size in each dimension is calculated
 358 as:

$$\text{Overlap} = K - S$$

359 The overlapping design reduces blocky artifacts and better preserves bound-
 360 aries and fine details because pixels near patch edges contribute to multiple
 361 patch vectors

Table 2: OPM parameters for hierarchical feature maps

Stage	Patch Size (K)	Stride (S)	Padding (P)	Output Resolution
1	7	4	3	$\frac{H}{4} \times \frac{W}{4}$
2	3	2	1	$\frac{H}{8} \times \frac{W}{8}$
3	3	2	1	$\frac{H}{16} \times \frac{W}{16}$
4	3	2	1	$\frac{H}{32} \times \frac{W}{32}$

362 The overlapping design improves segmentation performance because it sup-
 363 plies the transformer with smoother, more informative multi-scale features. Af-

364 ter patch merging, each stage’s feature maps are passed through transformer
 365 blocks, which include efficient self-attention and Mix-FFN layers.

366 4.1.3. Efficient Self-Attention (ESA)

367 SegFormer employs ESA, a computationally optimized adaptation of stan-
 368 dard self-attention used in ViTs. ESA is applied independently within each of
 369 the four stages of the MiT encoder. Given an input feature map $F_i \in \mathbb{R}^{H_i \times W_i \times C_i}$
 370 from stage i , ESA flatten spatial locations into a token sequence $X \in \mathbb{R}^{N \times C_i}$,
 371 where $N = H_i \cdot W_i$ represents the number of spatial locations. In standard
 372 multi-head self-attention the per-head queries, keys and values are computed as

$$Q = XW_Q, \quad K = XW_K, \quad V = XW_V, \quad (1)$$

373 where $W_Q, W_K, W_V \in \mathbb{R}^{C_i \times d_{\text{head}}}$ and d_{head} denotes the per-head dimension.
 374 Standard attention computes $\text{softmax}(\frac{QK^\top}{\sqrt{d_{\text{head}}}})V$, which requires forming an $N \times$
 375 N affinity matrix and therefore has quadratic complexity in the number of
 376 tokens.

377 ESA reduces this cost by shortening the key/value sequence by a reduc-
 378 tion ratio R . A sequence-reduction operator (denoted $\text{SeqReduce}(\cdot)$) produces
 379 downsampled keys and values

$$K' = \text{SeqReduce}(K) \in \mathbb{R}^{\frac{N}{R} \times d_{\text{head}}}, \quad V' = \text{SeqReduce}(V) \in \mathbb{R}^{\frac{N}{R} \times d_{\text{head}}}. \quad (2)$$

380 $\text{SeqReduce}(\cdot)$ can be implemented by reshaping and linear projection. ESA then
 381 computes attention from full-resolution queries to the reduced keys/values:

$$\text{EfficientAttention}(Q, K', V') = \text{softmax}\left(\frac{QK'^\top}{\sqrt{d_{\text{head}}}}\right)V', \quad (3)$$

382 where the softmax is taken along the reduced key dimension (length N/R) so
 383 that each of the N queries attends over the $\frac{N}{R}$ reduced positions. The com-
 384 putational cost becomes $\mathcal{O}(N \cdot \frac{N}{R} \cdot d_{\text{head}})$, which is significantly lower than the
 385 complexity of standard self-attention when $R > 1$ (Xie et al., 2021).

386 4.1.4. Mix Feed-Forward Network (Mix-FFN)

387 The Mix-FFN is a modification of the standard feed-forward network (FFN)
 388 that injects local spatial context into token-wise MLPs by inserting a 3×3

convolution between the two linear projections. This provides local positional information while preserving the global modelling capability of the FFN.

Let the input tokens be $x_{\text{in}} \in \mathbb{R}^{N \times C}$ with $N = H_i W_i$. The Mix-FFN proceeds as follows:

$$z = W_1 x_{\text{in}} + b_1 \in \mathbb{R}^{N \times d_{\text{exp}}}, \quad (4)$$

$$Z = \text{reshape}(z) \in \mathbb{R}^{H_i \times W_i \times d_{\text{exp}}}, \quad (5)$$

$$U = \text{Conv}_{3 \times 3}(Z; \text{padding} = 1) \in \mathbb{R}^{H_i \times W_i \times d_{\text{exp}}}, \quad (6)$$

$$V = \text{GELU}(U), \quad (7)$$

$$v = \text{flatten}(V) \in \mathbb{R}^{N \times d_{\text{exp}}}, \quad (8)$$

$$x_{\text{out}} = W_2 v + b_2 + x_{\text{in}} \in \mathbb{R}^{N \times C}. \quad (9)$$

where $W_1 : \mathbb{R}^C \rightarrow \mathbb{R}^{d_{\text{exp}}}$ and $W_2 : \mathbb{R}^{d_{\text{exp}}} \rightarrow \mathbb{R}^C$ are the two linear projections (MLPs) of the FFN and $d_{\text{exp}} = r \cdot C$ is the expansion dimension (commonly $r = 4$) (Xie et al., 2021). The 3×3 convolution uses padding 1 to preserve spatial resolution. The residual connection $+x_{\text{in}}$ is applied as in standard transformer blocks. After processing, the output is flattened back to $N \times C$ for subsequent layers.

4.2. Lightweight All-MLP Decoder

SegFormer’s decoder eliminates the complexity of traditional convolutional decoders by relying entirely on MLPs for efficient feature fusion and segmentation. The decoder unifies feature channel dimensions, upsamples features to a common spatial resolution, fuses them via a pointwise linear layer, and predicts per-pixel class logits with a final linear projection. For four-level encoder feature maps F_i the decoder proceeds as follows:

1. **Feature unification.** Each encoder feature map F_i (with C_i channels) is projected to a unified channel dimension C by a pointwise linear layer:

$$\hat{F}_i = \text{Linear}(C_i, C)(F_i) \quad \text{for all } i. \quad (10)$$

408 2. **Upsampling.** Each unified feature map \hat{F}_i is upsampled to the common
 409 spatial resolution $\frac{H}{4} \times \frac{W}{4}$:

$$\tilde{F}_i = \text{Upsample}\left(\frac{H}{4}, \frac{W}{4}\right)(\hat{F}_i) \quad \text{for all } i, \quad (11)$$

410 where \tilde{F}_i denotes the upsampled version of \hat{F}_i .

411 3. **Concatenation and fusion.** The upsampled features are concatenated
 412 along the channel dimension. For four encoder levels this yields $4C$ chan-
 413 nels, which are fused back to C channels by a pointwise linear layer:

$$F = \text{Linear}(4C, C)\left(\text{Concat}_i(\tilde{F}_i)\right). \quad (12)$$

414 4. **Segmentation prediction.** A final linear layer maps the fused feature
 415 F to per-pixel class logits for N_{cls} classes:

$$M = \text{Linear}(C, N_{\text{cls}})(F), \quad (13)$$

416 so that M has shape $\frac{H}{4} \times \frac{W}{4} \times N_{\text{cls}}$. M is typically upsampled (e.g.,
 417 bilinear) to the original image resolution $H \times W$ for evaluation and visu-
 418 alization.

419 4.3. Implementation description

420 Our framework uses the MiT-B3 backbone as the foundation. It was con-
 421 figured with four stages containing [3, 4, 18, 3] Transformer layers, respectively.
 422 The number of attention heads for the stages is [1, 2, 5, 8] and the corresponding
 423 embedding dimensions are [64, 128, 320, 512]. The lightweight all-MLP decoder
 424 employs a hidden dimension of 768 to fuse multi-scale features and produce
 425 segmentation outputs.

426 The model was trained for 3,000 iterations with a batch size of 4 using the
 427 AdamW optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\varepsilon = 10^{-8}$. The initial learning
 428 rates were set to $\eta_0^{\text{backbone}} = 6 \times 10^{-5}$ for the backbone and $\eta_0^{\text{decoder}} = 6 \times 10^{-4}$
 429 for the decoder. A polynomial learning-rate schedule was applied:

$$\eta_t = \eta_0 \left(1 - \frac{t}{T}\right)^{0.9},$$

where η_0 is the initial learning rate, t is the current iteration, and T is the total number of iterations.

The AGR-TR framework was implemented in PyTorch (v1.7.1) and trained on a Linux workstation equipped with two NVIDIA RTX A6000 GPUs (48 GB each). Model performance was evaluated on the original validation set to assess real-world applicability. For comparison, we implemented five baseline segmentation models using MMSegmentation (Contributors, 2020): DeepLabV3 (Chen et al., 2017), Segmenter (Strudel et al., 2021), K-Net (Zhang et al., 2021b) (K-Net), MaskFormer (Cheng et al., 2021), and U-Net (Ronneberger et al., 2015). All baselines were reimplemented within the same training and evaluation pipeline to ensure a fair comparison.

4.4. Evaluation metrics

The performance of the abnormal rice growth segmentation framework is evaluated using the Intersection over Union (IoU) metric (Wang et al., 2020), a standard measure for semantic segmentation that quantifies the pixel-wise overlap between predicted and ground-truth labels. For each class c we compute the pixel-level counts: true positives (TP_c), false positives (FP_c), and false negatives (FN_c). Here, TP_c is the number of pixels correctly predicted as class c , FP_c is the number of pixels incorrectly predicted as class c , and FN_c is the number of pixels belonging to class c but predicted as another class. The IoU for class c is defined as

$$IoU_c = \frac{TP_c}{TP_c + FP_c + FN_c}. \quad (14)$$

The mean IoU (mIoU) over N abnormality classes is computed as

$$mIoU = \frac{1}{N} \sum_{c=1}^N IoU_c, \quad (15)$$

where N denotes the total number of classes in the dataset. The mIoU penalizes both over- and under-segmentation and therefore provides a robust measure of segmentation accuracy.

455 To quantify uncertainty in the estimated performance, we report a 95%
 456 confidence interval (CI) for the mean mIoU computed across n independent
 457 experimental runs. Let x_i denote the mIoU observed in the i -th run, and define
 458 the sample mean and sample standard deviation by

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad (16)$$

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}. \quad (17)$$

459 Under the usual assumption that the sample mean is approximately t -distributed,
 460 a two-sided 95% CI for the true mIoU is given by

$$95\% \text{ CI} = \bar{x} \pm t_{\alpha/2, df} \cdot \frac{s}{\sqrt{n}}, \quad (18)$$

461 with $\alpha = 0.05$, degrees of freedom $df = n - 1$, and $t_{\alpha/2, df}$ the corresponding
 462 critical value from the Student’s t -distribution.

463 5. Experimental results

464 5.1. Data augmentation

465 To assess the effect of data augmentation on segmentation performance,
 466 each experiment (training with and without augmentation) was repeated for
 467 $n = 5$ independent runs using different fixed seeds. Table 3 summarizes ARG-
 468 TR’s segmentation performance on the original and augmented datasets. In
 469 the table “ \pm ” denotes the sample standard deviation across the n runs, and
 470 the 95% confidence intervals (CIs) for the mean mIoU were computed using the
 471 Student’s t -distribution with degrees of freedom $df = n - 1 = 4$.

472 The mean mIoU increased from 60.39% (original) to 62.88% (augmented),
 473 with corresponding 95% CIs [57.53%, 63.25%] and [60.15%, 65.61%], respec-
 474 tively. In addition, the consistent improvements on all evaluation metrics em-
 475 phasize the critical role of data augmentation in mitigating class imbalance
 476 challenges, refining feature learning, and enhancing generalization to diverse

Table 3: ARG-TR segmentation performance on original data and augmented data. Note: \pm indicates standard deviation.

	mIoU	mIoU 95% CI	Precision	Recall
Original data	60.39 \pm 2.3	[57.54, 63.24]	62.18 \pm 2.1	59.43 \pm 2.4
Augmented data	62.88 \pm 2.2	[60.15, 65.61]	65.03 \pm 2.7	63.82 \pm 2.9

field conditions. For example, the improvement in precision and recall suggests that the augmentation pipeline reduces both false positives and false negatives, where ambiguous or rare symptoms often challenge model robustness.

5.2. Spectral band contribution analysis

Figure 8 shows sample images for three different spectral-band configurations: green (G), green + near-infrared (G+NIR), and green + near-infrared + red-edge (G+NIR+RE).

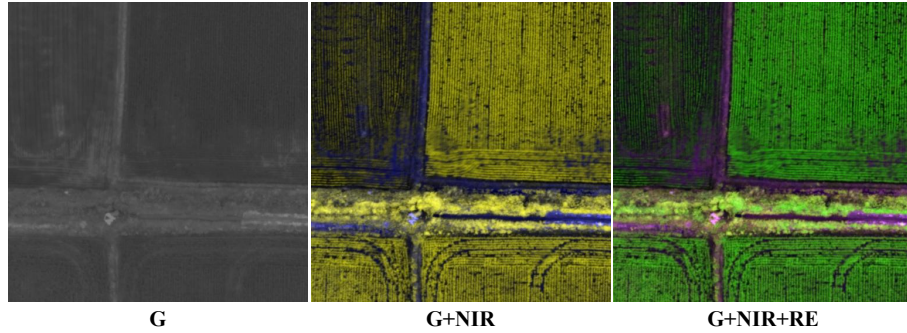


Figure 8: Illustration of the three spectral-band settings used as input to the model. “G” denotes the green band; “NIR” denotes near-infrared; “RE” represents the red-edge band.

To quantify the contribution of each spectral band to anomaly detection, we performed an ablation study using three input configurations: G only, G+NIR, and G+NIR+RE. Table 4 reports the class-wise IoU (%) for each configuration.

With only the green channel the model obtains moderate performance (IoU between 47.1% and 57.3%). The combination of G and NIR bands yield substantial gains for all anomaly types. For example, IoU for L increases from

Table 4: Ablation of spectral bands on class-wise IoU (%).

Class	G	G+NIR	G+NIR+RE
Missing plants (MP)	49.80	60.52	65.82
Lodging (L)	57.34	66.18	68.41
Rice blast (RBD)	47.11	57.73	61.78
Poor growth (PG)	50.52	58.21	63.34

57.3% to 66.1%, and IoU for MP increases from 49.8% to 60.5%. The integration of RE band further improves performance and produces the highest IoU for every class. For example, L to 68.4% and MP to 65.8%. These results indicate that NIR provides complementary contrast useful for detecting structural and vegetation anomalies, while the RE band refines discrimination of disease- and stress-related symptoms as it is sensitive to chlorophyll content and subtle stress signals. Overall, the combination G+NIR+RE offers the most informative spectral input for abnormal rice growth segmentation in our experiments.

5.3. ARG-TR performance evaluation

We performed an ablation study to examine how encoder depth and decoder hidden size affect segmentation accuracy and model complexity. Table 5 reports three ARG-TR variants with different encoder depths, decoder hidden dimensions, total parameter counts (in millions), and the resulting mIoU.

Table 5: Ablation of encoder depths and decoder hidden size. "Hidden sizes" lists the stage-wise embedding dimensions for the encoder.

Model variant	Depths	Hidden sizes	Decoder hidden size	Params (M)	mIoU
ARG-TR (1)	[2, 2, 2, 2]	[64, 128, 320, 512]	256	14.0	61.65
ARG-TR (2)	[3, 4, 6, 3]	""	768	25.4	62.72
ARG-TR (3)	[3, 4, 18, 3]	""	768	45.2	64.86

Key observations include:

- Moving from ARG-TR (1) to ARG-TR (2) increases the parameter count by 11.4M (from 14.0M to 25.4M, about 81.4% increase) and yields a smaller mIoU gain by 1.07% (from 61.65% to 62.72%).
- Moving from ARG-TR (2) to ARG-TR (3) further increases parameters by 19.8M (from 25.4M to 45.2M, approximately 78.0% increase) and obtains a larger mIoU value of 1.96% (from 62.72% to 64.86%).

These results show that increasing model capacity consistently improves segmentation performance, and in this set of variants the largest model (ARG-TR (3)) provides the highest mIoU. Considering the balance between accuracy and computational cost, ARG-TR (3) was selected as the primary model for subsequent experiments because it achieves the highest segmentation performance. For deployment scenarios with limited memory or latency budgets, ARG-TR (1) or ARG-TR (2) are preferable due to their lower parameter counts and competitive performance.

Figure 9 presents the training progress of the ARG-TR model using two key metrics recorded over 3,000 iterations: pixel accuracy (left) and training loss (right). The validation accuracy (seg_accuracy) shows a rapid rise between approximately 600 and 1,000 iterations, reaching about 88%–90%. The loss starts near 0.9 and decreases sharply to around 0.4 by iteration 1,500, then gradually stabilizes close to 0.4 by iteration 3,000. The quick initial convergence indicates that ARG-TR efficiently learns discriminative features even with limited labeled data. After the early-stopping mark at iteration 2,000, both accuracy and loss remain stable. Therefore, 2,000 iterations are considered sufficient to achieve near-optimal generalization in our setup.

Table 6 summarizes ARG-TR’s segmentation performance for four abnormal rice growth classes. Reported metrics are IoU, precision, and recall (all in percentage). The testing process was repeated for $n = 5$ independent runs for each class. Overall, ARG-TR achieves IoU over 60.8% for all classes. “Lodging” obtains the best performance (mean IoU = 67.3%, precision = 68.4%, recall = 69.13%), likely because its structural signature (bent or collapsed plants) is

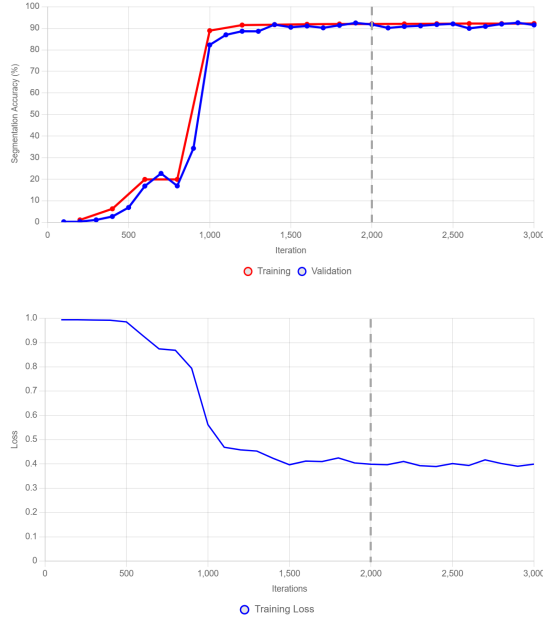


Figure 9: Training progress of the ARG-TR model on the abnormal rice growth dataset: pixel accuracy (left) and loss (right). The vertical dashed line indicates the early-stopping mark at iteration 2,000.

visually distinct. “Missing plants” follows with mean IoU = 64.1%.
 “Poor growth” and “Rice blast” show lower IoU values at 61.3% and 60.8%,
 respectively. The relatively lower precision and recall for “Poor growth” and
 “Rice blast” can be attributed to the higher visual similarity of these symptoms
 to healthy rice plants under certain conditions, which increases the likelihood
 of both false positives and false negatives. For “Poor growth”, the phenotypic
 differences, such as slight stunting, reduced leaf area, or lighter color, can be
 subtle and easily confused with natural field variability or early-stage nutrient
 deficiencies. For “Rice blast”, the appearance of lesions may be small for early-
 stage disease, sparsely distributed, or partially occluded by surrounding leaves.
 As a consequence, it is difficult to detect them at UAV imaging resolutions.

Table 6: ARG-TR performance for each abnormal rice growth class. Note: \pm indicates standard deviation.

	Missing plants	Poor growth	Lodging	Rice blast
IoU	64.11 \pm 2.2	61.29 \pm 2.5	67.37 \pm 1.8	60.87 \pm 2.8
Precision	66.80 \pm 2.5	64.27 \pm 2.1	68.42 \pm 2.5	63.47 \pm 2.3
Recall	67.34 \pm 1.8	63.98 \pm 2.4	69.13 \pm 2.4	62.18 \pm 2.0

5.4. Visualization of abnormal rice growth segmentation using the ARG-TR framework

Figure 10 and Figure 11 illustrate segmentation outputs produced by the ARG-TR framework for four abnormal rice growth classes: Missing Plants (MP), Poor Growth (PG), Rice Blast Disease (RBD), and Lodging (L). The color legend used throughout the figures is: MP (magenta), PG (red), RBD (cyan), L (orange), healthy rice (blue), and bare ground (black).

Overall, the visual alignment between model predictions and ground truth demonstrates robust segmentation performance for several categories. For MP, the model consistently detects large gaps in the field and closely matches ground truth boundaries. For L, the model successfully captures the irregular textures and patterns associated with lodged plants. For PG, the model locates small and sparse affected areas with relatively few false positives. Finally, the model correctly detects the infected RBD regions, which matches the ground truth. Figure 11 shows examples of mixed and ambiguous cases that highlight both strengths and weaknesses of the model.

In general, while the model effectively identifies large contiguous regions of L and RBD, it struggles with finer distinctions in overlapping or ambiguous cases. For example, MP areas are occasionally undersegmented or as L, particularly in regions where L regions are near the MP regions (Figure 11 c, d). Moreover, early stage RBD symptoms tend to be fragmented in predictions, which reflects the difficulty of separating infected regions from healthy areas. These errors highlight the complexity of identifying co-occurring stressors in real-world agri-

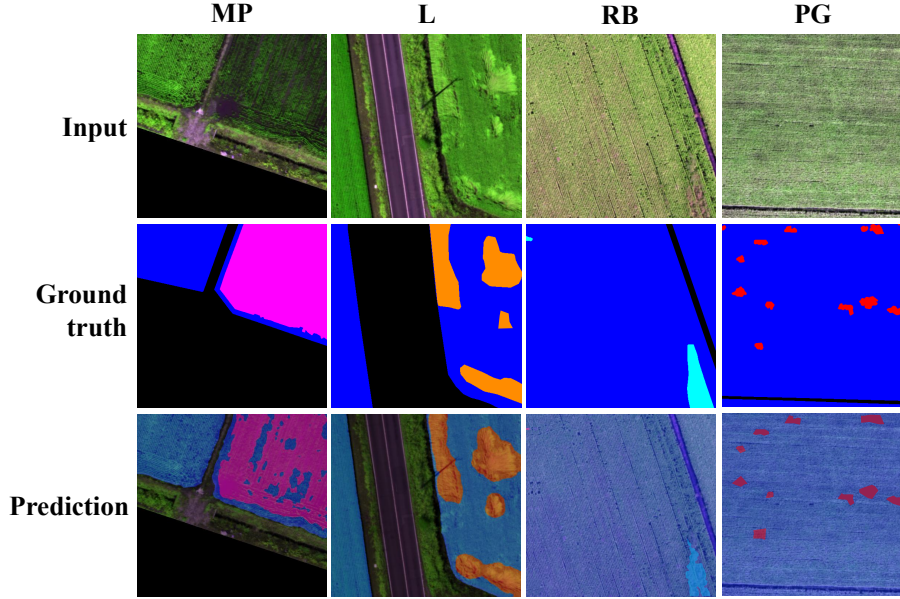


Figure 10: Predictions of the ARG-TR framework for all abnormal rice growth classes. **Note:** MP: Missing Plants (Magenta), PG: Poor Growth (Red), RBD: Rice Blast Disease (Cyan), L: Lodging (Orange), Blue indicates healthy rice, and Black represents bare ground.

568 cultural scenes, where symptom boundaries are often blurred by environmental
569 variability and plant interactions.

570 5.5. Comparative analysis of ARG-TR and other baseline segmentation models

571 This section compares the ARG-TR framework with several established seg-
572 mentation models: KNet (Zhang et al., 2021b), Segmenter (Strudel et al., 2021),
573 SegFormer (Xie et al., 2021), DeepLabv3 (Chen et al., 2017), U-Net (Ron-
574 neberger et al., 2015), MaskFormer (Cheng et al., 2021), and EDANet (Yang
575 et al., 2020). Table 7 summarizes each model’s performance on the abnormal rice
576 growth validation set using mIoU, pixel accuracy, and inference speed (frames
577 per second (FPS)).

578 ARG-TR achieves the highest mIoU (64.86%) and pixel accuracy (93.42%)

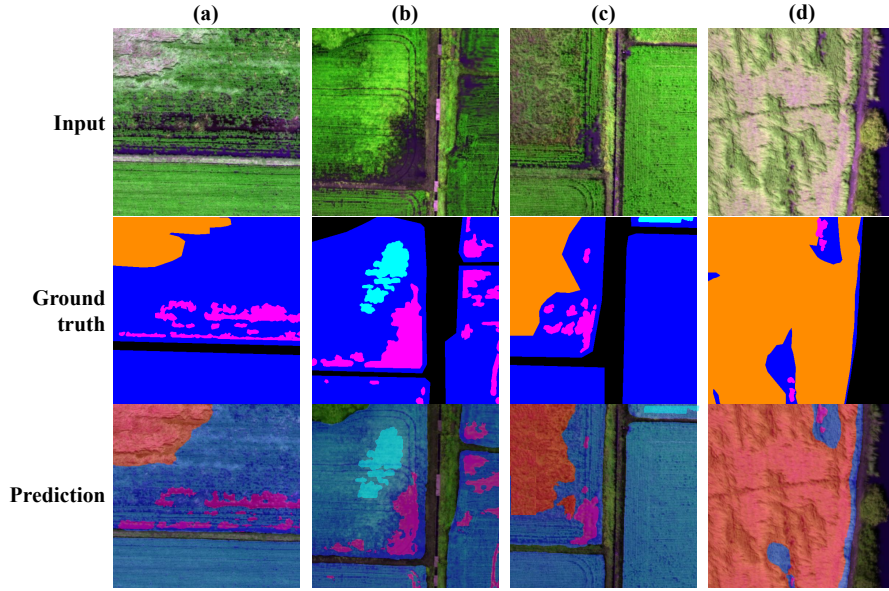


Figure 11: ARG-TR predictions for challenging mixed-condition cases. **Note:** MP: Missing Plants (Magenta), RBD: Rice Blast Disease (Cyan), L: Lodging (Orange), Blue indicates healthy rice, and Black represents bare ground.

579 on the dataset, outperforming strong baselines such as KNet (mIoU: 57.34%)
 580 and MaskFormer (mIoU: 60.13%). This improvement suggests that ARG-TR
 581 offers superior contextual understanding and finer feature discrimination, likely
 582 due to its transformer-based global-context modeling and the integration of
 583 targeted anomaly-aware modules.

584 Regarding efficiency, ARG-TR reaches a better trade-off between accuracy
 585 and speed. While it is slower than EDANet (61 FPS) and DeepLabv3 (29 FPS),
 586 its accuracy gains make it more suitable for precision agricultural monitoring
 587 where segmentation quality is prioritized. Lighter models, like UNet and some
 588 transformer variants, such as Segmenter and MaskFormer show lower segmenta-
 589 tion performance on this task, which highlights limitations in capturing complex
 590 spatial hierarchies.

591 Figure 12 presents qualitative comparisons between ARG-TR, MaskFormer,
 592 and KNet on three representative UAV samples. ARG-TR consistently produces

Table 7: Performance comparison of the ARG-TR framework and baseline models on the abnormal rice growth dataset. Note: Inference speed measured on the same evaluation environment (batch size = 1).

Model	mIoU	Pixel ac- curacy	Inference speed (FPS)
KNet (Zhang et al., 2021b)	57.34	89.12	12
Segmenter (Strudel et al., 2021)	56.47	88.75	10
DeepLabv3 (Chen et al., 2017)	54.89	86.43	29
UNet (Ronneberger et al., 2015)	49.21	83.56	15
MaskFormer (Cheng et al., 2021)	60.13	90.58	9
EDANet (Yang et al., 2020)	56.52	89.23	61
ARG-TR (Segformer) (Xie et al., 2021)	64.86	93.42	25

593 masks with sharp boundaries and reduced noise. In the first two samples (dis-
594 tinct L and MP regions), ARG-TR’s predictions align closely with ground truth
595 annotations. In the third sample, ARG-TR shows minor over-segmentation
596 but remains more similar to the ground truth than MaskFormer and KNet.
597 MaskFormer tends to produce more fragmented MP regions, while KNet pro-
598 duces noisier and more scattered masks, especially in samples with mixed ab-
599 normalities. These qualitative differences emphasize ARG-TR’s strengths in
600 fine-grained anomaly localization and boundary adherence, both important for
601 real-world agricultural monitoring where small or ambiguous symptoms must
602 be detected reliably.

603 Through previous experiments, ARG-TR consistently outperformed state-
604 of-the-art baselines in both mean IoU and pixel accuracy. Using G+NIR+RE
605 input bands produced a 13.2% increase in IoU compared with only using the
606 green channel. In addition, ARG-TR achieved real-time inference and produced
607 more precise segmentation boundaries, particularly in mixed or ambiguous field
608 conditions.

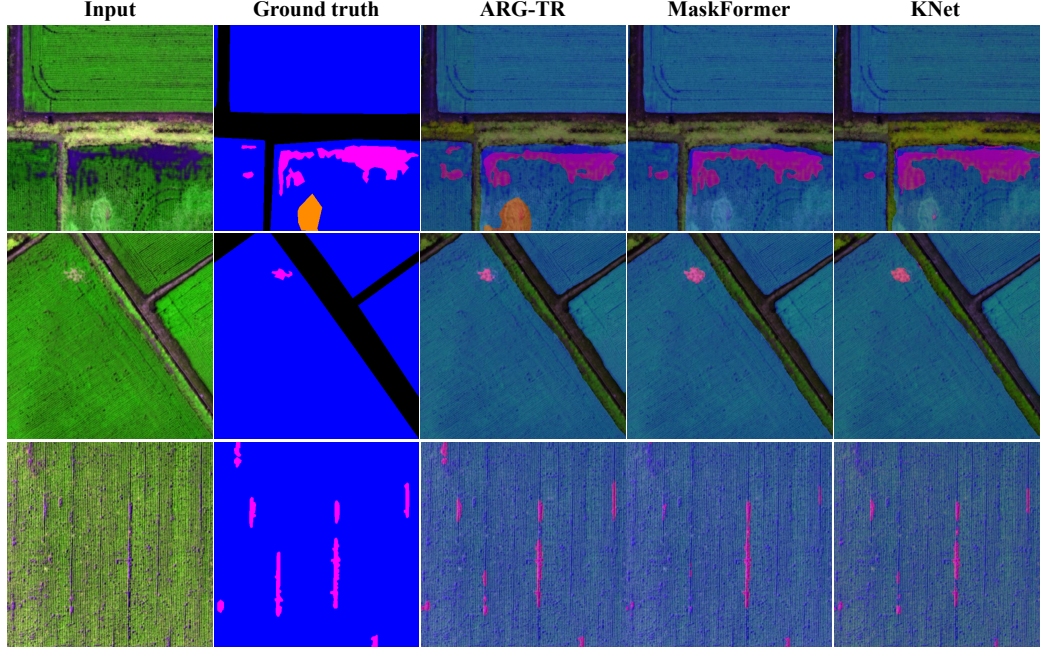


Figure 12: Comparison of the output of ARG-TR framework and two other state-of-the-art segmentation models, MaskFormer and KNet on three different input samples. **Note:** MP: Missing Plants (Magenta), L: Lodging (Orange), Blue indicates healthy rice, and Black represents bare ground.

6. Discussion

The primary goal of this study was to identify an efficient and robust DL framework for abnormal rice growth detection. We evaluated multiple segmentation architectures on a large, manually annotated UAV dataset. Through a series of experiments, transformer-based architectures, such as SegFormer and MaskFormer, achieved higher segmentation performance than other CNN-based alternatives (e.g., DeepLabv3, U-Net). These results are consistent with recent work that highlights transformers' ability to model global context and long-range dependencies for agricultural disease and stress detection (Wang et al., 2024; Kapetas et al., 2024). According to the results reported in Table 7, the ARG-TR framework showed the best overall performance ($mIoU = 64.86\%$,

620 pixel accuracy = 93.42%). The hierarchical feature fusion and transformer-
621 based global-context modeling improve discrimination of subtle anomalies such
622 as lodging and rice blast.

623 Several aspects of UAV data acquisition greatly affect the performance of
624 the framework. First, variation in solar illumination and viewing geometry in-
625 troduces spectral shifts that reduce class separability. We addressed this by
626 applying geometric correction, spectral band combination, and augmentations
627 during training, but residual effects can still increase false positives/negatives
628 in borderline cases. Second, weather constraints and the limited number of
629 imaging dates (five discrete growth stages) create temporal gaps that can miss
630 rapid symptom progression. For example, mIoU scores for Rice Blast Disease
631 (60.8%) and Poor Growth (61.3%) indicate lower per-class performance com-
632 pared with large, contiguous anomalies such as missing plants. This is expected
633 because stunting and early infections produce weak, spatially dispersed spec-
634 tral signatures that are difficult to distinguish from normal variability. Third,
635 flight altitude and ground-sampling distance limit the detectability of very small
636 or early-stage lesions; multispectral indices (e.g., NDVI, red-edge) partly com-
637 pensate by highlighting physiological stress that is not obvious in RGB, but
638 small-scale symptoms remain challenging. Finally, ARG-TR’s inference speed
639 (25 FPS measured in our evaluation setting) is suitable for many UAV-based
640 monitoring workflows where segmentation quality is prioritized. However, in ap-
641 plications that require very high throughput, such as continuous video streams
642 or large-area rapid surveys, lighter-weight models or optimized inference engines
643 (EDANet, DeepLabv3) are preferable.

644 Although this study focused on G, NIR, and RE bands, the network archi-
645 tecture can readily accommodate different spectral combinations or higher res-
646 olution sensors with minimal modification. While we demonstrated the model
647 mainly on rice, the framework can be retrained for other species, such as wheat
648 or maize, because the model learns spatial spectral representations directly from
649 data, it can adapt to diverse geographic regions and environmental conditions
650 given representative training samples. The model is suitable for real-time UAV

651 deployments or integration into monitoring platforms for various agricultural
652 settings.

653 **7. Conclusions and future works**

654 In this work, we introduced ARG-TR, a transformer-based segmentation
655 framework specifically, configured for identifying abnormal rice growth patterns
656 using drone-captured imagery. The model was trained on a large-scale drone-
657 based dataset containing 378,074 high-resolution images covering four common
658 abnormal rice growth anomalies (lodging, rice blast disease, poor growth, and
659 missing plants). By integrating hierarchical transformer architecture with a
660 strategic augmentation pipeline, ARG-TR achieves rapid convergence during
661 training and robust generalization to diverse field conditions. With a mIoU of
662 64.86 and 93.42% pixel accuracy, ARG-TR excels in identifying distinct anomalies
663 like lodging and rice blast disease, while maintaining efficient inference speed
664 (25 FPS).

665 Challenges exist in detecting subtle or overlapping stressors like early-stage
666 stunting and ambiguous symptom boundaries. Future work will explore hybrid
667 architectures that combine local texture encoders with global transformers, as
668 well as domain-specific synthetic augmentations to enrich rare-class representations.
669 Moreover, the integration of additional modalities, such as spectral
670 or temporal data, may further sharpen boundary delineation and symptom discrimination.
671 Finally, with continued enhancements in model design and training
672 strategies, ARG-TR has the potential to power real-time and scalable agriculture
673 systems capable of delivering timely and actionable insights for crop health
674 management.

675 **Acknowledgement**

676 This work was supported by the Young Scientists Fund of the National Natural
677 Science Foundation of China (No. 62502271), the Young Scientists Fund of
678 the Natural Science Foundation of Shandong Province (No. ZR2025QC630), the

679 Natural Science Foundation of Rizhao City (Nos. RZ2024ZR33 and RZ2024ZR34)
 680 and by Institute of Information & communications Technology Planning & Eval-
 681 uation (IITP) under the metaverse support program to nurture the best talents
 682 (IITP-2024-RS-2023-00254529) grant funded by the Korea government(MSIT).

683 **Author contributions**

684 **Yanfen Li**: Writing – original draft. **L. Minh Dang**: Methodology, Writ-
 685 ing – original draft, Writing – review & editing, Data curation. **Hanxiang**
 686 **Wang**: Conceptualization, Writing – review & editing. **Muhammad Fayaz**:
 687 Data curation. **Sufyan Danish**: Visualization. **Junliang Shang**: Formal anal-
 688 ysis, Supervision. **Hyoung-Kyu Song**: Funding acquisition. **Hyeonjoon**
 689 **Moon**: Supervision.

690 **Declaration of Competing Interest**

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

694 **References**

- 695 Agüera-Vega, F., Carvajal-Ramírez, F., & Martínez-Carricondo, P. (2017). As-
 696 sessment of photogrammetric mapping accuracy based on variation ground
 697 control points number using unmanned aerial vehicle. *Measurement*, 98,
 698 221–227.
- 699 Alam, N., Sagar, A. S., Dang, L. M., Zhang, W., Park, H. Y., & Hyeonjoon, M.
 700 (2025). Deep learning based radish and leaf segmentation for phenotype trait
 701 measurement. *Signal, Image and Video Processing*, 19, 178.
- 702 Bin Rahman, A. R., & Zhang, J. (2023). Trends in rice research: 2030 and
 703 beyond. *Food and Energy Security*, 12, e390.

- 704 Biswal, S., Pathak, N., Chatterjee, C., & Mailapalli, D. R. (2024). Estimation of
705 aboveground biomass from spectral and textural characteristics of paddy crop
706 using uav-multispectral images and machine learning techniques. *Geocarto*
707 *International*, *39*, 2364725.
- 708 Chen, L.-C., Papandreou, G., Kokkinos, I., Murphy, K., & Yuille, A. L. (2017).
709 Deeplab: Semantic image segmentation with deep convolutional nets, atrous
710 convolution, and fully connected crfs. *IEEE transactions on pattern analysis*
711 *and machine intelligence*, *40*, 834–848.
- 712 Cheng, B., Schwing, A., & Kirillov, A. (2021). Per-pixel classification is not all
713 you need for semantic segmentation. *Advances in neural information process-*
714 *ing systems*, *34*, 17864–17875.
- 715 Contributors, M. (2020). Mmsegmentation: Openmmlab semantic segmentation
716 toolbox and benchmark.
- 717 Dang, L. M., Hassan, S. I., Suhyeon, I., kumar Sangaiah, A., Mehmood, I., Rho,
718 S., Seo, S., & Moon, H. (2020). Uav based wilt detection system via convo-
719 lutional neural networks. *Sustainable Computing: Informatics and Systems*,
720 *28*, 100250.
- 721 Dang, M., Wang, H., Li, Y., Nguyen, T.-H., Tightiz, L., Xuan-Mung, N., &
722 Nguyen, T. N. (2024). Computer vision for plant disease recognition: a com-
723 prehensive review. *The Botanical Review*, *90*, 251–311.
- 724 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Un-
725 terthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S. et al. (2020).
726 An image is worth 16x16 words: Transformers for image recognition at scale.
727 *arXiv preprint arXiv:2010.11929*, .
- 728 Dosset, A., Dang, L. M., Alharbi, F., Habib, S., Alam, N., Park, H. Y., &
729 Moon, H. (2025). Cassava disease detection using a lightweight modified soft
730 attention network. *Pest Management Science*, *81*, 607–617.

Food, & of the United Nations, A. O. (2025). Food and agriculture projections to 2050. <https://www.fao.org/global-perspectives-studies/food-agriculture-projections-to-2050/en/>, accessed 2025-04-03.

Ju, O.-J., Choi, B.-R., Jang, E. K., Soh, H., Lee, S.-W., & Lee, Y.-S. (2022). Climate change and rice yield in hwaseong-si gyeonggi-do over the past 20 years (2001~ 2020). *Korean Journal of Environmental Agriculture*, 41, 16–23.

Kang, Y., Meng, Q., Liu, M., Zou, Y., & Wang, X. (2021). Crop classification based on red edge features analysis of gf-6 wfv data. *Sensors*, 21, 4328.

Kapetas, D., Kalogeropoulou, E., Christakakis, P., Klaridopoulos, C., & Pechlivani, E. M. (2024). Multi-spectral image transformer descriptor classification combined with molecular tools for early detection of tomato grey mould. *Smart Agricultural Technology*, 9, 100580.

Li, Y., Wang, H., Dang, L. M., Sadeghi-Niaraki, A., & Moon, H. (2020). Crop pest recognition in natural scenes using convolutional neural networks. *Computers and Electronics in Agriculture*, 169, 105174.

Lin, T., Wang, Y., Liu, X., & Qiu, X. (2022). A survey of transformers. *AI open*, 3, 111–132.

Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., & Guo, B. (2021). Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision* (pp. 10012–10022).

Measurusa (2025). TRINITY F90+. <https://measurusa.com/pages/trinity-f90>, accessed 2025-04-05.

MicaSense (2025). MicaSense RedEdge-MX. <https://support.micasense.com/hc/en-us/articles/360011389334-RedEdge-MX-Integration-Guide>, accessed 2025-04-05.

758 Nuradili, P., Zhou, J., & Melgani, F. (2024). Wetland segmentation method
759 for uav multispectral remote sensing images based on segformer. In *IGARSS*
760 *2024-2024 IEEE International Geoscience and Remote Sensing Symposium*
761 (pp. 6576–6579). IEEE.

762 Rezvi, H. U. A., Tahjib-Ul-Arif, M., Azim, M. A., Tumpa, T. A., Tipu, M.
763 M. H., Najnine, F., Dawood, M. F., Skalicky, M., & Brestič, M. (2023). Rice
764 and food security: Climate change implications and the future prospects for
765 nutritional security. *Food and Energy Security*, *12*, e430.

766 Ronneberger, O., Fischer, P., & Brox, T. (2015). U-net: Convolutional networks
767 for biomedical image segmentation. In *International Conference on Medical*
768 *image computing and computer-assisted intervention* (pp. 234–241). Springer.

769 Samal, P., Babu, S. C., Mondal, B., & Mishra, S. N. (2022). The global rice
770 agriculture towards 2050: An inter-continental perspective. *Outlook on Agri-*
771 *culture*, *51*, 164–172.

772 Spasev, V., Dimitrovski, I., Chorbev, I., & Kitanovski, I. (2024). Semantic seg-
773 mentation of unmanned aerial vehicle remote sensing images using segformer.
774 In *International Conference on Intelligent Systems and Pattern Recognition*
775 (pp. 108–122). Springer.

776 Stankus, A. (2021). State of world aquaculture 2020 and regional reviews: Fao
777 webinar series. *FAO aquaculture newsletter*, (pp. 17–18).

778 Strudel, R., Garcia, R., Laptev, I., & Schmid, C. (2021). Segmenter: Trans-
779 former for semantic segmentation. In *Proceedings of the IEEE/CVF interna-*
780 *tional conference on computer vision* (pp. 7262–7272).

781 Tian, M., Ban, S., Yuan, T., Ji, Y., Ma, C., & Li, L. (2021). Assessing rice
782 lodging using uav visible and multispectral image. *International journal of*
783 *remote sensing*, *42*, 8840–8857.

784 Wang, H., Nguyen, T.-H., Nguyen, T. N., & Dang, M. (2024). Pd-tr: End-to-
785 end plant diseases detection using a transformer. *Computers and Electronics*
786 *in Agriculture*, 224, 109123.

787 Wang, Z., Wang, E., & Zhu, Y. (2020). Image segmentation evaluation: a survey
788 of methods. *Artificial Intelligence Review*, 53, 5637–5674.

789 Wu, S., Ma, X., Jin, Y., Yang, J., Zhang, W., Zhang, H., Wang, H., Chen,
790 Y., Lin, C., & Qi, L. (2025). A novel method for detecting missing seedlings
791 based on uav images and rice transplanter operation information. *Computers*
792 *and Electronics in Agriculture*, 229, 109789.

793 Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J. M., & Luo, P. (2021).
794 Segformer: Simple and efficient design for semantic segmentation with trans-
795 formers. *Advances in neural information processing systems*, 34, 12077–
796 12090.

797 Yang, M.-D., Boubin, J. G., Tsai, H. P., Tseng, H.-H., Hsu, Y.-C., & Stewart,
798 C. C. (2020). Adaptive autonomous uav scouting for rice lodging assessment
799 using edge computing with deep learning edanet. *Computers and Electronics*
800 *in Agriculture*, 179, 105817.

801 Zhang, T., Xu, Z., Su, J., Yang, Z., Liu, C., Chen, W.-H., & Li, J. (2021a). Ir-
802 unet: Irregular segmentation u-shape network for wheat yellow rust detection
803 by uav multispectral imagery. *Remote Sensing*, 13, 3892.

804 Zhang, W., Pang, J., Chen, K., & Loy, C. C. (2021b). K-net: Towards unified
805 image segmentation. *Advances in Neural Information Processing Systems*,
806 34, 10326–10338.