

NortheastChinaSoybeanYield20m: an annual soybean yield dataset at 20 m in Northeast China from 2019 to 2023

Jingyuan Xu^{1,2}, Xin Du^{1,2}, Taifeng Dong³, Qiangzi Li^{1,2}, Yuan Zhang^{1,2}, Hongyan Wang^{1,2}, Jing Xiao^{1,2},
5 Jiashu Zhang^{1,4}, Yunqi Shen^{1,2}, Yong Dong^{1,2}

¹Aerospace Information Research Institute, Chinese Academy of Sciences, Beijing 100094, China

²University of Chinese Academy of Sciences, Beijing 100190, China

³National Wildlife Research Centre, Environment and Climate Change Canada, 1125 Colonel By Drive, Ottawa, ON K1A0H3, Canada

10 ⁴School of Science, China University of Geosciences (Beijing), Beijing 100083, China

Correspondence to: Xin Du (duxin@aircas.ac.cn)

Abstract. Accurate monitoring of crop yield is critical for ensuring food security. While various yield datasets covering Northeast China exist, they were produced at a coarse spatial resolution and remain inadequate for capturing small-scale spatial heterogeneity. Current yield estimation methods, such as machine learning models and the assimilation of remotely sensed
15 biophysical variables into crop growth models, are heavily reliant on ground observations and computationally expensive. To address these limitations, we propose a hybrid framework that couples the World Food Studies Simulation Model (WOFOST) and a Gated Recurrent Unit (GRU) model to generate a high-resolution (20 m) soybean yield dataset in Northeast China from 2019 to 2023 (NortheastChinaSoybeanYield20m). First, to generate a comprehensive training dataset, WOFOST was employed to simulate diverse soybean growth scenarios by accounting for variations in climates, crop varieties, soil types and
20 agro-managements practices. The GRU model was then trained to establish relationships between model simulated leaf area index (LAI) and soybean yield. The trained model was applied to estimate soybean yield in Northeast China using time-series LAI derived from Sentinel-2 at key growth stages. The accuracy of estimates was evaluated using in-situ measurements and government statistical data. The overall accuracy was 287.44 kg ha⁻¹ and 272.36 kg ha⁻¹ in the root mean squared error (RMSE) for field and regional scale, respectively. The model exhibited consistent interannual stability, with mean relative error (MRE)
25 averaging 11.46 % and 7.94% at the municipal scale and the provincial scale, respectively. The dataset effectively captured spatiotemporal yield variability, offering potentials for optimizing soybean production, guiding precise agriculture practices, and informing agricultural policy. The NortheastChinaSoybeanYield20m dataset is publicly available at <https://doi.org/10.5281/zenodo.14263103> (Xu et al., 2024).

1 Introduction

30 Soybean is a crucial crop for both food and oil production, providing more than a quarter of the world's edible protein (Graham and Vance, 2003). Global demand for soybean is projected to increase by 46 % by 2050, driven by rapid population growth (Falcon et al., 2022). As an major traded agricultural commodity, soybean production in key exporting nations has wide-reaching effects on international markets, and can significantly influence agricultural economies worldwide (Qiao et al., 2023). Notably, China is the world's largest consumer of soybeans (FAOSTAT, 2022), and its soybean demand relies heavily on
35 international trade (Zhao et al., 2023). Consequently, accurate monitoring of soybean yield is vital for promoting sustainable agriculture, ensuring food security, and maintaining economic stability from regional to global scale. Moreover, effective yield monitoring and mapping supports farmers by informing field management practices, bolstering agricultural insurance and enhancing poverty alleviation initiatives (Zhuo et al., 2022).

Remote sensing data provides time-series observations for crop yield estimation across multiple scales (e.g., field,
40 regional and national) (Dong et al., 2020; Hunt et al., 2019; Zhao et al., 2023b). Current methodologies for yield estimation can be broadly categorized as data-driven or knowledge-driven approaches.

Data-driven methods leverage satellite-derived variables such as leaf area index (LAI), fraction of absorbed photosynthetically active radiation (FAPAR), and vegetation indices (VIs) to establish linear or nonlinear relationships with measured crop yield (Ang et al., 2022; Xie et al., 2019). Machine learning algorithms such as Random Forest (RF), and
45 Artificial Neural Networks (ANN), due to their ability to process large dataset and model complex nonlinear interactions, have been widely applied in crop yield estimations (Pang et al., 2022; Tian et al., 2021; Yildirim et al., 2022). These methods can extract effective information from multi-source structured or unstructured data without manual intervention. However, they are heavily reliant on extensive ground-truth training data, which is challenging to collect over large areas and high time intervals (Cao et al., 2021). Additionally, these models often overlook the impacts of environmental factors on crop growth,
50 such as the influence of early-season soil moisture on root establishment or the effect of high temperatures during flowering on pod set, and are lack of interpretability, as they cannot explain the causal relationship between input features and outputs, leading to poor spatial-temporal generalization (Gevaert, 2022).

In contrast, knowledge-driven crop growth models simulate crop development from sowing to harvest based on agronomic mechanisms (Kaur and Singh, 2020). Common model types include light-use efficiency models (e.g., SAFY
55 (Duchemin et al., 2008)), soil-driven models (e.g., AquaCrop (Steduto et al., 2009)), and atmospheric-driven models (e.g., WOFOST (Diepen et al., 1989)). These models integrate environmental factors (e.g., climate conditions and soil characteristics) with crop physiological processes (Gasol et al., 2024). Climate variables like temperature, precipitation, and solar radiation are critical in regulating essential physiological processes such as photosynthesis, respiration and transpiration, which influence the rate and duration of crop growth stages (Misaal et al., 2023). Climate anomalies during specific growth stages may disrupt
60 biochemical processes, ultimately affecting yield formation. Similarly, soil properties influence crop productivity by regulating water retention, aeration, and nutrient uptake (Muhuri et al., 2023). Despite their mechanistic rigor, applications of crop models

over large area are typically constrained by (1) insufficient spatial-temporal input data, and (2) parameter uncertainty, which can propagate errors into yield estimations (Dokoohaki et al., 2021). To overcome these challenges, data assimilation techniques to integrate remote sensing observations (e.g., LAI) into crop growth models have been developed to enhance spatial representativity (Huang et al., 2024). However, high resolution remote sensing data drastically increases computational cost, limiting the scalability of these approaches for regional or national mappings efforts (Huang et al., 2019).

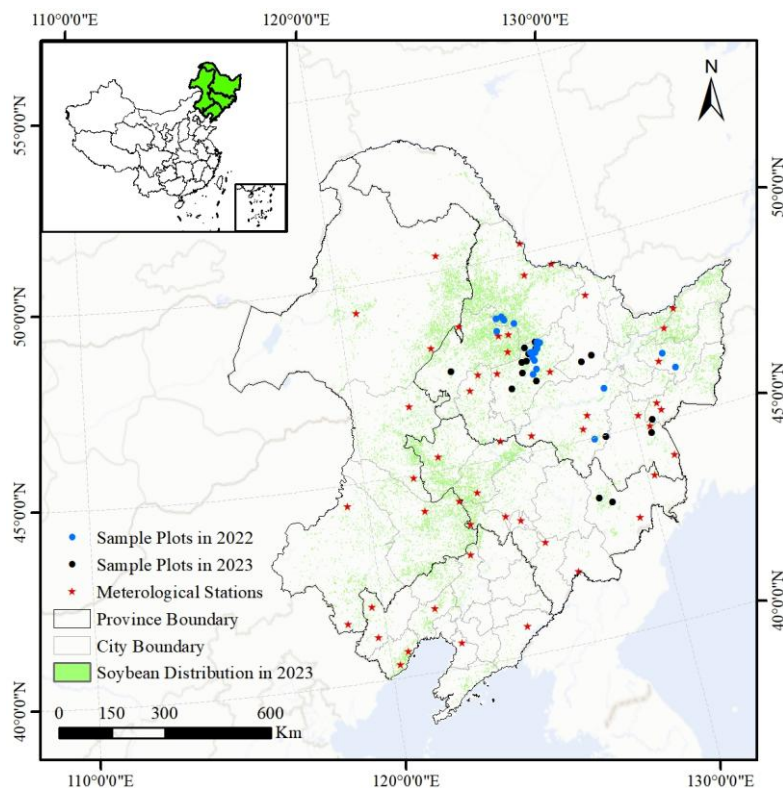
Given the limitations above, integrating data-driven and knowledge-driven models has emerged as a critical strategy to enhance spatial-temporal generalization and mitigate sparse training data challenges in crop yield estimations. Hybrid frameworks coupling crop growth model with machine learning algorithm, such as those proposed and evaluated by Ren et al., (2023b) and Xie and Huang, (2021), are gaining tractions. These approaches utilized simulated outputs from crop growth models (e.g., meteorological, soil, crop physiological, and management factors) as inputs for machine learning, reducing reliance on limited ground observations. Many studies have demonstrated hybrid methods are able to enhance yield estimation due to three benefits (Feng et al., 2020; Xie and Huang, 2021; Yang et al., 2021). The simulations from crop growth model can provide biophysical constraints to machine learning, ensuring agronomic plausibility. The crop growth models generate synthetic training datasets to address data scarcity. Finally, the machine learning improves the computational efficiency compared to traditional data assimilation techniques (Xie and Huang, 2021). However, existing studies generally extracted input features (e.g., LAI, and soil moisture) across the entire growth cycle or on coarse temporal scales, increasing computational costs of model calculation and obscuring stage-specific physiological response (Pinke and Lövei, 2017; Wang et al., 2015). Additionally, while deep learning models, such as Long Short-Term Memory (LSTM) and GRU model excel at modelling temporal dependencies, their integration into hybrid frameworks have not been widely explored.

Critically, the primary soybean-producing regions of China lack a publicly available high-resolution yield dataset to analyse spatiotemporal production patterns, hindering precision agriculture and policy optimization. To address this, we developed a hybrid model coupling the World Food Studies (WOFOST) crop growth model with a GRU deep learning method to estimate soybean yield in Northeast China. The objectives include: (1) Design a hybrid framework integrating WOFOST-simulated growth scenarios with GRU-based temporal feature extraction; (2) Generate a high-resolution (20 m) soybean yield dataset in Northeast China (NortheastChinaSoybeanYield20m) from 2019 to 2023; (3) Evaluate the accuracy of the dataset across field, municipal, and provincial scales using in situ and statistical benchmarks. The WOFOST model first simulated a multi-scenario soybean growth (varying climate, soil, crop varieties and management conditions) to train the GRU model. The time series Sentinel-2 data, capturing soybean growth development, were then input into the GRU model to estimate yield. This approach prioritizes stage-specific physiological dynamics which balancing computational efficiency and spatial granularity, providing a critical advancement for scalable agricultural monitoring.

2 Data preparation and preprocessing

2.1 Study areas

The study was conducted in Northeast China (38°40' N to 53°34' N, 115°05' E to 135°02' E), encompassing Heilongjiang, Jilin, Liaoning province, as well as the eastern parts of the Inner Mongolia Autonomous Region (IMAR) (Fig. 1). The study area includes 40 cities and spans approximately 1.24 million km². The region is characterized by a continental monsoon climate, with an annual accumulated temperature ($\geq 10\text{ }^{\circ}\text{C}$) ranging from 2200 to 3600 $^{\circ}\text{C}$ (Pu et al., 2019), and a frost-free period of 140 to 170 days (Tan et al., 2014). The average annual precipitation exhibits a strong east-west gradient, decreasing from 1000 mm in the east to 350 mm in the west (Zhao et al., 2011). The predominant soil types include brown coniferous forest soil, dark brown forest soil, forest steppe chernozem and meadow grassland chernozem soil (Pu et al., 2019). Soybean is one of the three main crops in the region, primarily cultivated in the northern parts of the Songliao plain in rotation with maize. Notably, this region contributes around 64 % of China's total annual soybean production (National Bureau of Statistics of China (NBSC), 2023). Approximately 97 % of the soybean in the region is rainfed (Guo et al., 2022; Yu et al., 2020), with growing season typically spanning from May to late September (Zhao et al., 2021).



105

Figure 1: Location of the study area and the distribution of sample plots in two years (2022 and 2023) and selected meteorological stations. The soybean distribution map was obtained from Zhao et al., (2022) using a moment-preserving segmentation method, achieving an overall accuracy over 90% for soybean in 2023 (Details are provided in Section 2.2.5).

2.2 Data collections

110 2.2.1 In-situ measurement data

Due to limitations of resources and personnel, in-situ measurements were not available during the earlier years (from 2019 to 2021). Field-scale yield data was separately collected through field investigation in September 2022 and 2023. In each year, a total of 21 and 18 sample plots were selected, respectively (Fig. 1). Within each sample plot that was around 100 m × 100 m in area, nine quadrats with area of 1 m × 1 m were selected randomly for destructive sampling of yield in soybean. The central
115 location of each quadrat was recorded using a GPS device with accuracy of 1 m. The harvested beans were then oven-dried about 72 hours in Hailun Agricultural Ecology Experimental Station, Chinese Academy of Sciences to determine the yield. Finally, the average yield for the selected nine quadrats represents the soybean yield of the sample plot. In addition, soybean planting dates for different regions were collected through field surveys, providing agro-management data for this study.

Field measured LAI data of soybean was obtained from the Common Application Support Platform for Land Observation
120 Satellite (CAPLOS, <https://124.16.188.131:9699/web/server3/build/#/Guide>), an open data portal providing in-situ biophysical variables (e.g., LAI and vegetation cover) for validating remote sensing products and refining retrieval algorithms. LAI measurements were collected using a LICOR LAI-2200 plant canopy analyzer following a standardized protocol. At each site, the instrument was positioned above the canopy to obtain a reference reading of incoming solar radiation, and then positioned about 5 cm above ground to collect six readings of radiation transmitting beneath the soybean canopy. The raw data
125 was taken rigorous quality control to remove outliers, missing or duplicate values. After processing, 94 LAI observations were retained, spanning three soybean growing seasons (2021 – 2023).

2.2.2 Meteorological data

In this study, two different climate datasets were used.

The meteorological station data used in this study came from the meteorological stations of the National Meteorological
130 Information Center (<http://data.cma.cn>). There are 238 meteorological stations within the study area. Here 51 of the meteorological stations that located within 1 km buffer zone of the soybean cultivation areas were selected (Fig. 1). The meteorological datasets generally include insolation duration (h), minimum temperature (°C), maximum temperature (°C), daily average temperature (°C), average water vapor pressure (kPa), average wind speed (m sec⁻¹), precipitation (mm) and snow-depth (cm). Observed data from 1980 to 2021 of the 51 selected stations were collected. Missing values and outliers in
135 the data were filtered out. The data were then directly used for setting input climate parameters of the WOFOST model to drive simulations.

The climate reanalysis data was obtained from the ERA5-land Daily Aggregated - ECMWF Climate Reanalysis Product. The data was only used to calculate soybean phenology for preparation of yield estimations. It was a global climate reanalysis product that provides continuous climate data at a resolution of 0.1° × 0.1° (e.g., air temperature and atmospheric pressure)
140 starting from 1950. The daily aggregated air temperature data at 2 m above the surface of land measured in kelvin (K) during

the soybean growth periods from 2019 to 2023 was collected in this study from the Google Earth Engine (<http://earthengine.google.com>). The product was resampled to 20 m using bilinear interpolation model to match with the resolution of satellite imagery data.

2.2.3 Soil data

145 Soil data was obtained from the 1:1000,000 Chinese soil database, established by the Institute of Soil Science, Chinese Academy of Sciences (Shi et al., 2004). The dataset consisted of two parts: soil spatial data (digital soil maps) and soil attribute data. In this study, the 1:1000,000 soil spatial data was obtained. The spatial database was developed by digitizing, mosaicking, and reassembling sheets from the 1:1,000,000 Soil Map of the People's Republic of China (National Soil Survey Office, 1995), with the Genetic Soil Classification of China (GSCC) soil families as the fundamental mapping units. The final dataset includes
150 909 soil types and over 94,000 polygons. The dataset was utilized to determine the dominant soil types within the study area, serving as the basis for assigning soil parameter settings according to literatures.

2.2.4 Satellite imagery data

Two satellite data including: 1) Sentinel-2 Multi-Spectral Instrument (MSI) Level - 2A Surface reflectance product (10 – 60 m spatial resolution, 5-day revisit), and 2) the Moderate Resolution Imaging Spectroradiometer (MODIS) Leaf Area Index (LAI) / Fraction of Photosynthetically Active Radiation (FPAR) Level 4 product (MCD15A3H, v061, 500 m spatial resolution,
155 4-day period) were used to generate yield maps. All data spanning soybean growth periods (2019 – 2023) were accessed and pre-processed via the Google Earth Engine (GEE, <http://earthengine.google.com>).

The MSI aboard Sentinel-2A/B satellites provides 10 m (visible and near-infrared bands), 20 m (red-edge and shortwave infrared bands) and 60 m (atmospheric bands) bands at 5-day revisit. The Level-2A data, which are geometrically and
160 atmospherically corrected via the Sen2Cor, were masked for clouds and shadows using the Quality Assurance (QA) band. The 60 m band was excluded due to their low spatial resolution and limited relevance for yield estimation and the 10 m (B2: Blue, B3: Green, B4: Red, B8: Near-Infrared) and 20 m (B5–B7: Red-edge, B8A: Near-Infrared, B11–B12: Shortwave Infrared) bands were retained. To harmonize spatial resolution, the 10 m bands were resampled to 20 m using bilinear interpolation model.

165 The MODIS MCD15A3H (Collection 6.1, Level 4) provides 4-day composite LAI and FAPAR at 500 m derived from Terra and Aqua satellite sensors LAI/FAPAR are primarily inverted via a 3D radiative transfer model-based look-up-table (LUT) algorithm (Knyazikhin et al., 2018). When the primary algorithm fails, they are estimated using an empirical NDVI-LAI model. The LAI data was similarly reprojected to WGS -84 to ensure spatial alignment with Sentinel-2 imagery. These coarse-resolution LAI data were used to generate 500 m yield maps. The coarse-resolution yield maps were then used to bias-
170 correct the 20 m Sentinel-2 yield maps, improving their regional consistency. Details about the bias correction are present in following 3.3.2 Section.

2.2.5 Crop distribution data

175 The soybean distribution maps for the study area (2019 – 2023) were obtained from Zhao et al., (2022), which employed a novel methodology for crop type identification. The study proposed an optimal identification feature (OIF) knowledge graph coupled with a moment-preserving segmentation method to classify crop types without ground-truth data. The method achieved overall accuracy above 90% and producer’s accuracy exceeding 93% for maize, soybean and rice, with a Kappa coefficient greater than 0.90.

2.2.6 Statistical data

180 Crop yield records (1980-2022) were obtained from the Statistical Yearbooks published by the Statistic Bureau of Heilongjiang (<http://tjj.hlj.gov.cn>), Jilin (<http://tjj.jl.gov.cn>), Liaoning (<https://tjj.ln.gov.cn>) and Inner Mongolia Autonomous Region (<https://tj.nmg.gov.cn>) to validate the crop yield estimates. Because the 2022 Statistical Yearbook was not fully released, yield records for that year cover only a subset of cities. The statistical data served two main purposes, model simulation validation and regional-scale accuracy evaluation in this study. To ensure the multi-scenario soybean growth dataset capture the full range of production conditions that across multi-years meteorological data, various soil types, multiple soybean varieties and different agro-managements, the yield records from 1980 to 2022 along with published yield data and field samples were used to assess the reasonableness of simulated yields. For the spatial validation, regionally aggregated statistical yield data (2019 – 2022) were applied to evaluate the accuracy of the hybrid framework at municipal and provincial scales.

3 Methodology

190 Our proposed hybrid model utilizes both the advantages of machine learning in data mining and the mechanism advantages of crop growth model. Figure 2 presents the flowchart of the hybrid methodology for soybean yield estimation. It mainly includes 1) Generating a training dataset based on the WOFOST model that simulate multi-scenario soybean growth and yields under various climates, soil, cultivars and agro-management practices, 2) Training a GRU model to identify the relationships between simulated LAI and yield, 3) Producing soybean yield maps under multi-scale using LAI derived from MODIS and Sentinel-2 remote sensing data.

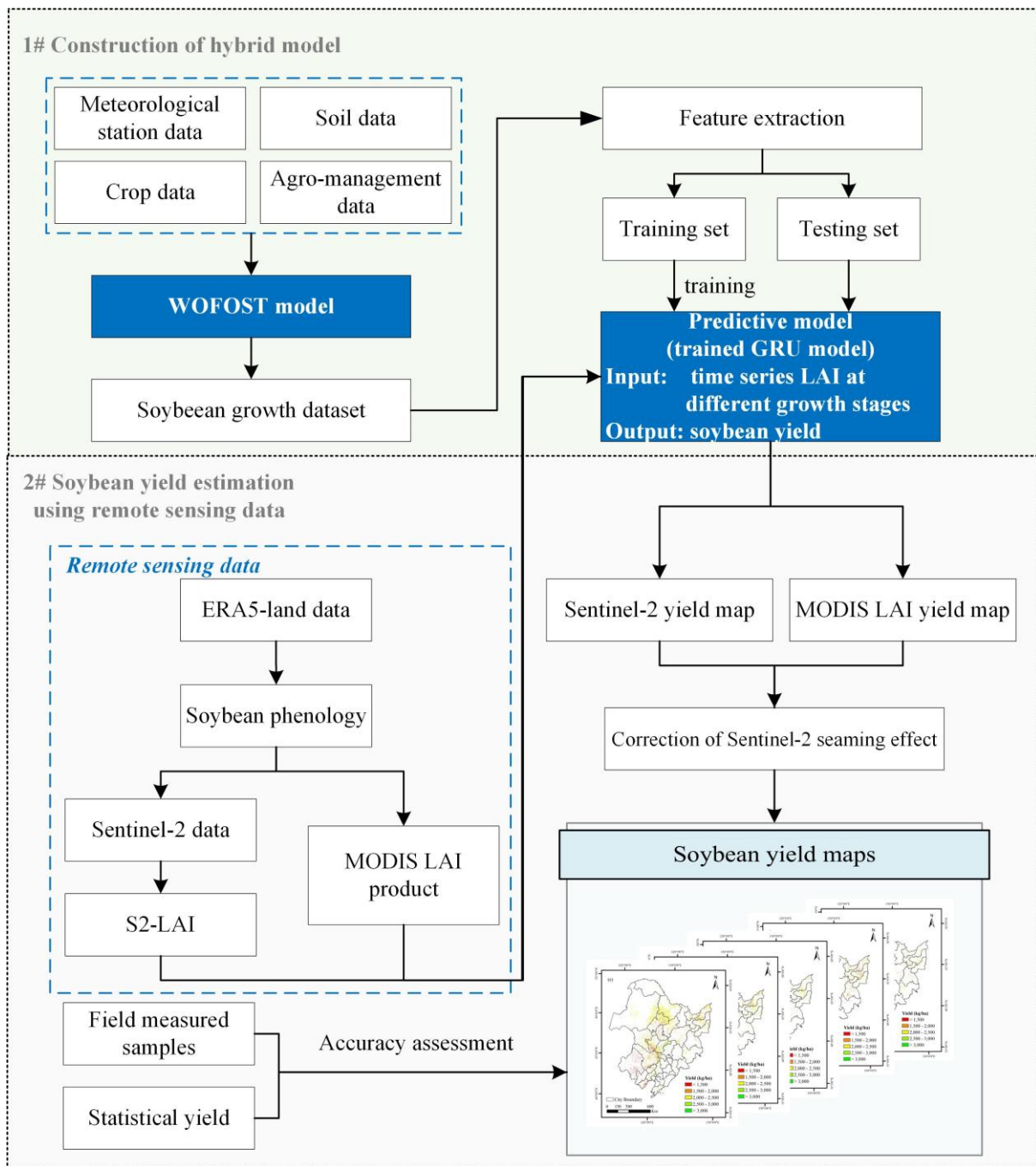


Figure 2: The flowchart of the overall yield estimation methodology in this study.

3.1 Construction of multi-scenario soybean growth dataset

The soybean growth dataset used in this study was a multi-scenario, knowledge-based dataset derived from crop growth model simulations. The dataset provided a quantitative description of yield formation by simulating soybean development under diverse agricultural production scenarios including variations in meteorology (temperature, precipitation, solar radiation), soil types (texture, organic matter), crop varieties (phenology, thermal requirements) and agro-managements (sowing date).

To generate the dataset, we employed the World Food Studies Simulation Model (WOFOST) (Diepen et al., 1989), implemented via the Python Crop Simulation Environment framework (PCSE, v5.5). The WOFOST model is well-suited for large-scale simulations and has been extensively validated (Huang et al., 2015). Given that soybean cultivation in the study region is predominantly rainfed, we adopted the water-limited mode (Wofost72_WLP_CWB) for simulations. Driven by daily weather data, soil profiles, and cultivar parameters, the water-limited model simulated CO₂-driven photosynthesis, water deficits, and biomass partitioning, outputting daily LAI, daily biomass accumulation, and final grain yield. Crop development is modelled through development stages (DVS) : 0 for emergence, 1 for anthesis and 2 for maturity (Diepen et al., 1989). This dataset mechanistically captures yield-limiting processes (e.g., drought during critical growth phases) while enabling scalable scenario analysis across Northeast China’s rainfed soybean systems.

3.1.1 Preparation of model input parameters

In this study, parameters of WOFOST were not calibrated or optimized using in-situ data, as our goal was to generate a scenario-based dataset representing diverse agricultural conditions rather than global optimal value. The model parameters were sourced from: ground observation, online database, peer-reviewed literatures and default WOFOST values.

(1) Meteorological parameters

The meteorological parameters required in WOFOST is shown in Table 1. To capture regional climate variability (e.g., temperature extremes, rainfall patterns), meteorological data of the selected 51 meteorological stations spanning 42 years (1980-2021) were compiled. These data – including daily temperature, precipitation, and solar radiation – were preprocessed into the model’s required input format (e.g., daily time steps, unit conversions) to ensure compatibility.

Table 1 Meteorological parameters required in WOFOST.

Parameter	Description	Units
IRRAD	Incoming global shortwave radiation	KJ m ⁻² d ⁻¹
TMIN	Daily minimum temperature	°C
TMAX	Daily maximum temperature	°C
VAP	Daily mean vapour pressure	kPa
WIND	Daily mean windspeed at 2 m above the surface	m s ⁻¹

RAIN	Daily rainfall	mm
SNOWDEPTH	Snow depth	cm

(2) Soil parameters

The soil parameters in the WOFOST mainly include soil moisture content at wilting point (SMW), field capacity (SMFCF) and saturation (SM0) as well as hydraulic conductivity of saturated soil (K0). Based on the 1:1,000,000 Chinese soil database, the study area predominantly comprises loam soil that is further classified into sandy, light, medium and heavy loam. The parameters for sandy, loam and medium loam were sourced from Du et al., (2025), while the parameters for heavy loam came from Sun et al., (2022). All soil parameter values, summarized in Table 2, were integrated into the model to evaluate the influence of soil variability on soybean yield (Du et al., 2025; Sun et al., 2022).

Table 2 Values of main soil parameters in WOFOST.

Soil type	SMW (cm ³ cm ⁻³)	SMFCF (cm ³ cm ⁻³)	SM0 (cm ³ cm ⁻³)	K0 (cm d ⁻¹)
Sandy loam	0.060	0.280	0.350	22.6
Light loam	0.090	0.280	0.340	19.3
Medium loam	0.110	0.280	0.340	18.1
Heavy loam	0.194	0.355	0.356	34.6

(3) Crop-specific parameters

In this study, the soybeans were classified into five types including early, medium-early, intermediate, medium-late and late maturity according to Qu et al., (2023). In the WOFOST model, soybean phenology is governed by temperature-driven parameters: the minimum (TBASEM) and maximum (TEFFMX) threshold temperature for emergence, and accumulated thermal time (TSUMEM: sowing to emergence; TSUM1: emergence to anthesis; TSUM2: anthesis to maturity). These thermal parameters are cultivar-sensitive and were set based on historical meteorological data and field phenology records, validated against field observations (Qu et al., 2023). Remaining crop parameters (e.g., SLATB: specific leaf area) were assigned default values or optimal values from Sun et al., (2022). Full parameter specifications are provided in Table A1.

(4) Agro-management parameters

Planting date is the major agro-management factors for soybean in the study area. The difference of planting date can significantly impact on soybean growth development, pod count, and biomass accumulation (Urda et al., 2024). Four planting dates 20 April, 30 April, 10 May, and 20 May to reflect the typical sowing window (late April to late May) of the study area were set for model simulation according to Mei et al., (2024).

3.1.2 Multi-scenarios crop simulations

Following parameter preparation, the four parameter categories, including meteorological (51 stations × 42 years), soil (4 types), crop-specific (5 varieties) and agro-management (4 planting dates), were systematically combined to create 171,360

unique scenarios (Table 3). These scenarios were executed in the WOFOST simulations, yielding a dataset of 171,360 various simulations that quantify yield responses to diverse agricultural production conditions.

Table 3 Scenarios for WOFOST simulations

Parameters	Number of categories	Details
Meteorological parameters	51×42	Meteorological data from 51 stations over 42 years (1980 – 2021)
Soil parameters	4	Sandy loam, light loam, medium loam and heavy loam
Crop-specific parameters	5	Early maturity, medium-early maturity, intermediate maturity, medium-late maturity and late maturity
Agro-management parameters	4	Four planting dates 20 April, 30 April, 10 May, and 20 May

3.2 Development of the Grated Recurrent Unit model (GRU)

250 A GRU (Grated Recurrent Unit) model, a streamlined variant of recurrent neural networks (RNNs), was employed to be trained using the multi- scenarios simulated dataset for large-scale soybean yield estimation. Unlike LSTM (Long short-term memory), GRU simplifies gating mechanisms to two adaptive gates, update and reset gates (Cho et al., 2014). The update gate retains the past information for future calculations. The reset gate aims to remove irrelevant historical context for simplifying the new candidate hidden states. Using the two gates together is benefit to balance long-term dependency capture and computational efficiency (Peng and Yili, 2022; Zhang et al., 2022). This design mitigates vanishing gradient issues while accelerating model training, making GRU particularly effective for time-series yield estimation (Gopi and Karthikeyan, 2023; Ren et al., 2023b).

260 Trained on the multi-scenarios simulated dataset, the GRU constructed based on TensorFlow 2.6 linked simulated environmental inputs to yield outputs. Accounting for the computational efficiency of the model in large areas, two key features include LAI_{mean1} (mean LAI during vegetative growth: emergence to flowering) and LAI_{mean2} (mean LAI during reproductive growth: flowering to maturity), were calculated to reflect photosynthetic capacity and yield potential. These two LAI metrics served as inputs, while simulated yields acted as outputs. The multi-scenarios simulated dataset was partitioned using 10-fold cross-validation, with hyperparameters (e.g. learning rate and batch size) optimized using a grid search to achieve minimal root mean squared error (RMSE, Eq. (5)) (Açikkar, 2024).

Once trained, the GRU model taken Sentinel-2-derived LAI time series as inputs to generate 20 m yield maps.

3.3.1 Determination of soybean phenology

Soybean phenology of the study area exhibits significant spatial variability due to climatic and varietal differences (Gaso et al., 2024). To address this, soybean phenology maps were generated from daily thermal time by integrating thermal zone divisions and regional adapted cultivars (Fig. A1). Soybean phenology (including emergence, anthesis and maturity) were calculated using daily aggregated air temperature data from ERA5-land dataset and a thermal time model (T_e) (Eq. (1)):

$$T_e = \begin{cases} 0, & (T_{mean} \leq T_{base}) \\ T_{mean} - T_{base}, & (T_{base} < T_{mean} < T_{max}) \\ T_{max} - T_{base}, & (T_{mean} \geq T_{max}) \end{cases} \quad (1)$$

where T_{mean} is daily mean temperature, T_{base} (8 °C) and T_{max} (37 °C) represent the minimum and the maximum temperature for soybean development, respectively (Allen et al., 1997; Choi et al., 2016). Soybean growth proceeds from a growth stage to the next stage when accumulated T_e reached the threshold of accumulated temperature required for growth according to the setting of crop parameters of WOFOST model (Table A1).

Based on field surveys and literatures, the planting dates of soybean were fixed as 5 May for Heilongjiang Province and Inner Mongolia Autonomous Region, and 1 May for Jilin and Liaoning Province (Huang and Liu, 2024), with emergence constrained to before 1 June (Mei et al., 2024), and maturity to before 1 October (Huang and Liu, 2024). Under the constraints of T_e (Table A1) and the agro-management, phenological dates (emergence, anthesis and maturity) were computed for each Sentinel-2 pixel (2019 – 2023). Finally, the phenological maps were clustered into 10 phenology classes using K-means clustering method (Jain and Dubes, 1988), aligning with Sentinel-2's revisit cycle to optimize imagery selection for yield estimation.

3.3.2 Model estimations of soybean yield

The red-edge normalized difference vegetation index ($NDVI_{RE}$) (Gitelson and Merzlyak, 1994) was employed for LAI mapping (Eq. (2)).

$$NDVI_{RE} = \frac{NIR-RE}{NIR+RE} \quad (2)$$

where B8A (near-infrared) and B5 (red-edge) are Sentinel-2 bands.

The soybean LAI was estimated from $NDVI_{RE}$ using a linear regression (Eq. (3)) validated across multiple crops ($R^2 = 0.732$, RMSE = 0.69) (Pasqualotto et al., 2019).

$$LAI = 5.405 \cdot NDVI_{RE} - 0.114 \quad (3)$$

Mean LAI values for vegetative (LAI_{mean1}) and reproductive (LAI_{mean2}) growth stages were computed from the time-series S2-derived LAI. These LAI values were input into the GRU model for yield prediction. For pixels with missing values during these stages, LAI values were replaced by the average of eight neighbouring pixels. Final 20 m yield maps were marked using soybean distribution maps to exclude non-soybean areas.

295 For large area estimations, a total of 194 Sentinel-2 tiles were required to fully cover the study area. Affected by cloud cover, the frequency of available data varied across each tile. Therefore, the yield maps often exhibited discontinuities along the edges of different tiles (“seaming effects”). This seaming effect could obscure real yield variations. To address this issue, a bias correction method proposed by Azzari et al., (2017) was applied. The overall framework is to use yield estimation based on MODIS LAI to correct the yield estimation from Sentinel-2. MCD15A3H generally provided more continuous estimation results of LAI due to its higher temporal resolution (4-day composites) and broader coverage. Yield maps were generated from the trained GRU taking MCD15A3LAI products as inputs. Sentinel-2 yield maps were adjusted by adding the difference between MODIS-derived mean yield and initial Sentinel-2 mean yield for each tile. This process minimized seams while preserving fine-scale yield variability within tiles.

3.4 Accuracy evaluation

305 The accuracy of generated NortheastChinaSoybeanYield20m (2019-2023) was evaluated on multiple scales. For field scale, in-situ yield data in 2022 and 2023 was used for assessment. For regional scale, the mean soybean yield for each city and province were separately calculated for each year, and compared with the statistical data. Accuracy evaluation was based on the coefficient of determination (R^2 , Eq. (4)), the root mean squared error (RMSE, Eq. (5)) and mean relative error (MRE, Eq. (6)).

$$310 \quad R^2 = 1 - \frac{\sum_i (y_{o,i} - y_{m,i})^2}{\sum_i (y_{o,i} - \bar{y}_0)^2} \quad (4)$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_{o,i} - y_{m,i})^2}{n}} \quad (5)$$

$$MRE = \frac{\sum_{i=1}^n |y_{o,i} - y_{m,i}|}{n \cdot y_{o,i}} \quad (6)$$

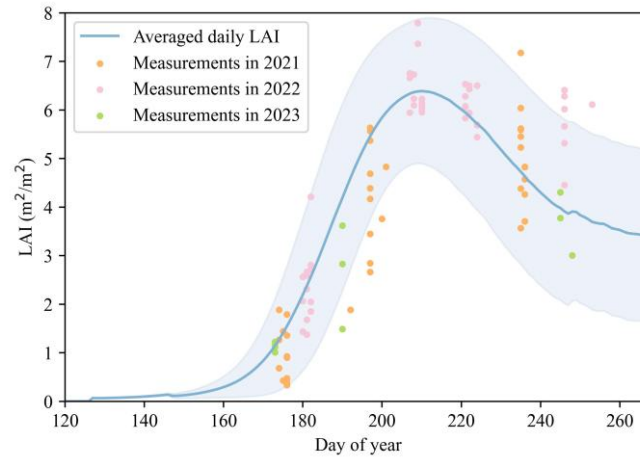
where $y_{o,i}$ and $y_{m,i}$ represent the actual yield (observed or statistical yield) and model estimated yield, respectively, \bar{y}_0 is the mean value of the actual yield.

315 4 Results and analysis

4.1 Simulations of the WOFOST model

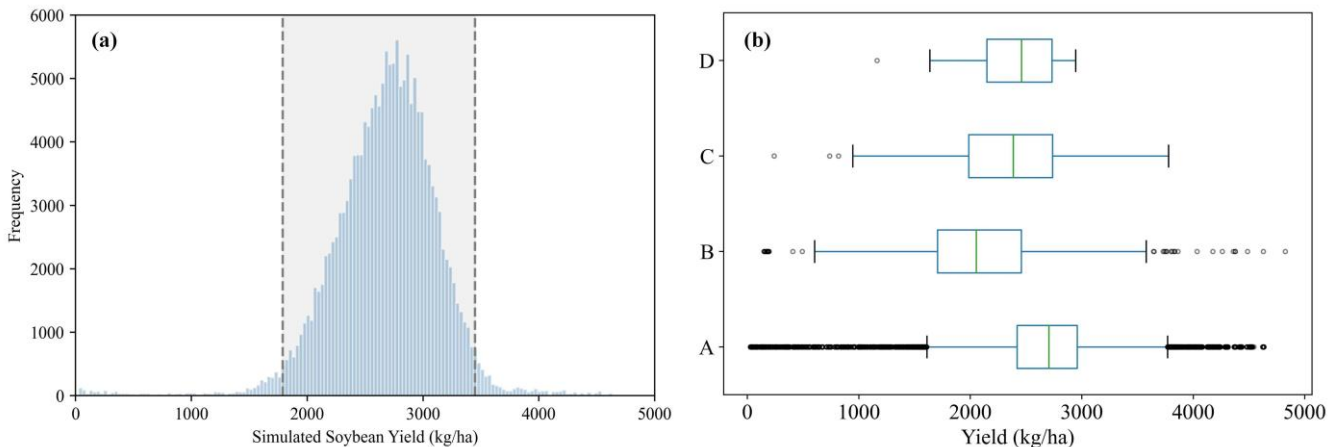
Since LAI was used as the input feature, the accuracy of WOFST-simulated LAI directly influenced the reliability of the GRU model for soybean yield prediction. To validate simulated LAI, field-measured LAI from 2021 to 2023 were compared against mean LAI curve calculated from value of 5,000 simulated LAI curves randomly selected from the multi-scenarios simulated dataset (Fig. 3). The results showed that simulated LAI trends aligned closely with observed field variations, capturing 88 % of the field-measured sample sites ($n = 83$) within the simulated range. This demonstrated robust agreement

between model outputs and ground truth, confirming the WOFOST simulations' ability to represent realistic LAI dynamics for GRU training.



325 **Figure 3: Comparison of averaged daily LAI randomly selected from model simulations (n = 5000) with field-measured LAI in 2021 (n = 38), 2022 (n = 46) and 2023 (n = 10). The gray shading represents one standard deviation, indicating the uncertainty in LAI simulation.**

Figure 4 (a) displays the histogram distribution of simulated soybean yields, revealing an approximately normal distribution (mean = 2675.66 kg ha⁻¹). The multi-scenario soybean growth dataset effectively captured wide range of production conditions, spanning low to high yield extremes. Fig. 4 (b) shows a box plot for comparing the simulated yield with historical statistics from 1980 to 2022, published yield data from literatures, and field measurements from 2022 and 2023. The simulated dataset exhibited the widest value range, demonstrating the comprehensiveness of the multi-scenario knowledge base and the robustness of the simulation outcomes.

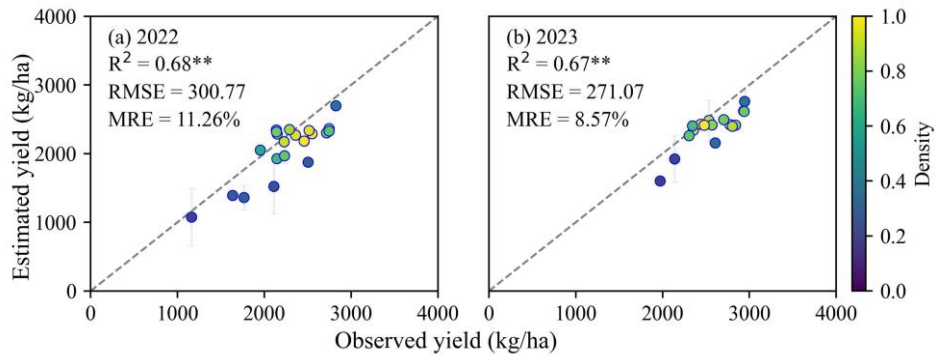


335 **Figure 4: (a) Histogram statistics of simulated soybean yield where the gray area in the histogram represents 95 % confidence intervals; (b) distribution of simulated soybean yield compared with other datasets where A represents simulated yield in this study (n = 171,360), B represents statistical yield from 1980 to 2022 (n = 961), C represents specific measurements from the literature**

(Chen et al., 2011; Fan et al., 2012; Liu et al., 2005, 2008; Liu and Herbert, 2002; Wang et al., 2020, 2024; Zheng and Zhang, 2021) (n = 138) and D represents measurements in 2022 and 2023 carried by this study (n = 39).

340 4.2 Yield estimation at field scale

The field-scale performance of NortheastChinaSoybeanYield20m was validated against in-situ measurement from 2022 and 2023, demonstrating strong accuracy in capturing spatial yield variability (Fig. 5). The estimated yields showed strong agreement with observed yield, with $R^2 > 0.65$ ($p < 0.01$). The error-bars indicated more consistent performance in fields with uniform yields, while higher uncertainties appear in fields with larger estimation deviations. Overall accuracy across both years reached 0.73 in R^2 ($p < 0.01$), 287.44 kg ha⁻¹ in RMSE and 10.02 % in MRE (Fig. A2). Notably, higher accuracy in 2023 with RMSE of 271.07 kg ha⁻¹ and MRE of 8.57 % (Fig. 5b) was achieved. The results indicated that the dataset well captured the spatial variation of soybean yield.

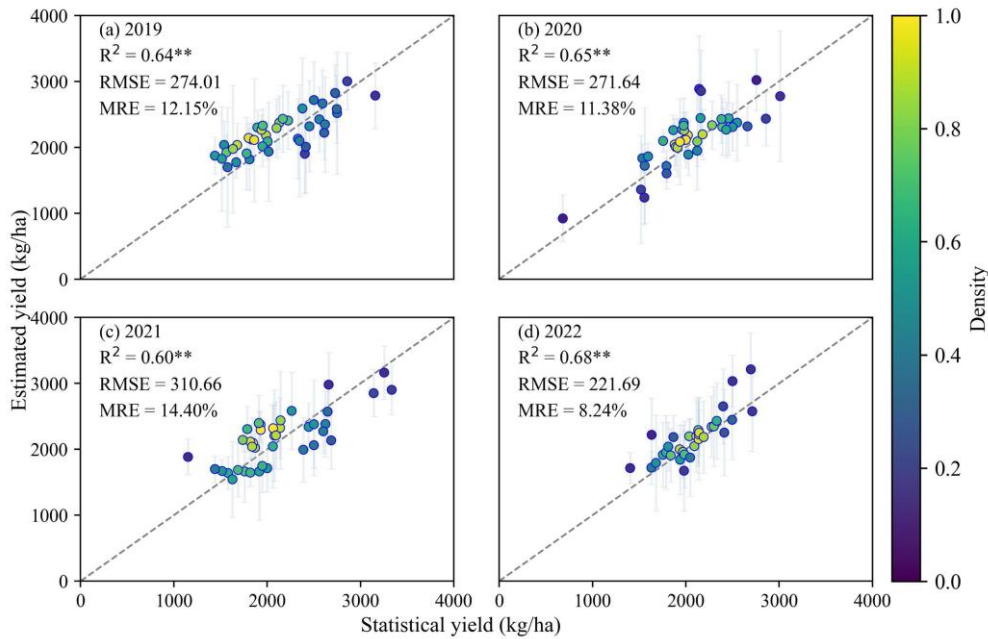


350 **Figure 5: Scatterplots between estimated and observed soybean yield in 2022 and 2023, respectively. The error-bars represent one standard deviation indicating the uncertainty of yield estimations. Dashed line represents 1:1 line. ** denotes statistical significance at $p < 0.01$.**

4.3 Yield estimation at regional scale

4.3.1 Variability of accuracy through years

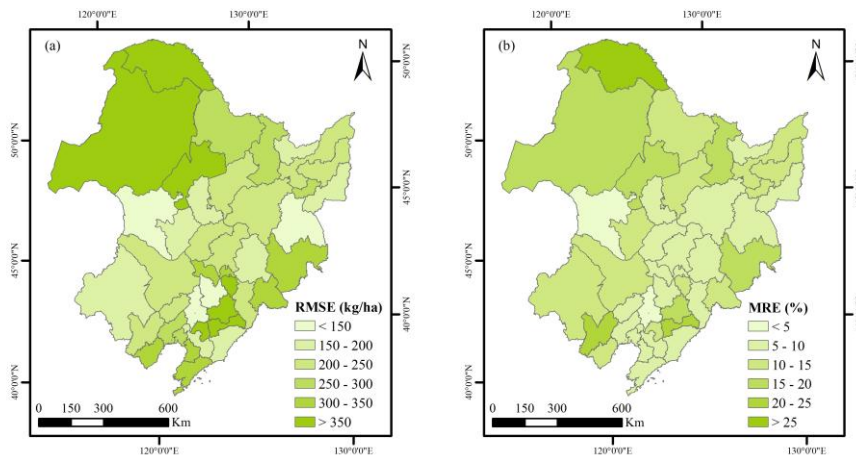
The NortheastChinaSoybeanYield20m was validated at the municipal scale (2019 to 2022) by aggregating yield maps to match statistical data (Fig. 6). Compared to the field-scale validation, the municipal-scale estimates exhibited greater uncertainty, likely reflecting increased heterogeneity of soybean yields over larger areas. The estimates maintained stable interannual performance, with correlation between estimated and statistical yields consistently exceeding 0.60 ($p < 0.01$). The overall accuracy, pooled across 2019- 2022, for municipal-scale achieved $R^2 = 0.62$ ($p < 0.01$), RMSE = 272.36 kg ha⁻¹, and MRE = 12.08 % (Fig. 11a). Annual accuracy metrics ranged from 221.69 kg ha⁻¹ to 310.66 kg ha⁻¹ for RMSE and from 8.24 % to 14.40 % for MRE, with the 2022 year achieving the highest accuracy (MRE < 10%, Fig. 6d).



365 **Figure 6: Scatterplots between estimated soybean yield from Sentinel-2 and municipal statistical yields for 2019 – 2022 (excluding 2023 for which no statistical data was not available from the government). The error-bars represent one standard deviation indicating the uncertainty of yield estimations. Dashed line represents 1:1 line. ** denotes statistical significance at $p < 0.01$.**

For temporal performance at the municipal scale, the RMSE between estimated and statistical yields from 2019 to 2022 remained below 500 kg ha^{-1} , with 80 % of cities exhibiting RMSE under 350 kg ha^{-1} (Fig. 7a). Spatially, large errors were concentrated in the northern part of Northeast China especially for the Greater Khingan Mountains area, while the flatter region, central regions of Northeast China, showed smaller deviations. The spatial distribution pattern of MRE closely mirrored that of RMSE (Fig. 7b), averaging value of 11.46 % across all cities over the four-year period. These findings underscore the

370 model's robust capacity to capture interannual variability of soybean yield.



375

Figure 7: Spatial patterns of the mean value of the root mean squared error and mean relative error between model estimated yields from Sentinel-2 and statistical yields from 2019 to 2022 (excluding 2023 for which no statistical data was not available from the government), (a) and (b), respectively. For years from 2019 to 2021, a total of 40 cities were calculated. For 2022, 32 cities were calculated due to missing statistics.

4.3.2 Spatial-temporal dynamics of soybean yield

380

To examine spatial patterns of soybean yield across Northeast China, yield distribution maps for 2019 - 2023 were generated (Fig. 8a-e). After bias correction with estimated yield derived from MODIS LAI products, the Sentinel-2 striping artifacts were markedly reduced (Fig. 8 vs. uncorrected estimates in Fig. A3), resulting in seamless 20 m yield surfaces with strong spatial continuity. Detailed yield estimations can be found in Fig. 9. Across five-year estimates, soybean yield in Northeast China predominantly ranged between 1500 and 2500 kg ha⁻¹, with higher yield concentrated in the central part plains where a region characterized by flat terrain and factorable agroclimatic conditions. The predicted yield is consistent with the municipal-scale statistical data (Table 4). Spatial variability, quantified by the coefficient of variation (CV), ranged from 17.51 % to 29.65 % over the study period, reflecting both inter and intra-annual heterogeneity in soybean productivity (Table 4).

385

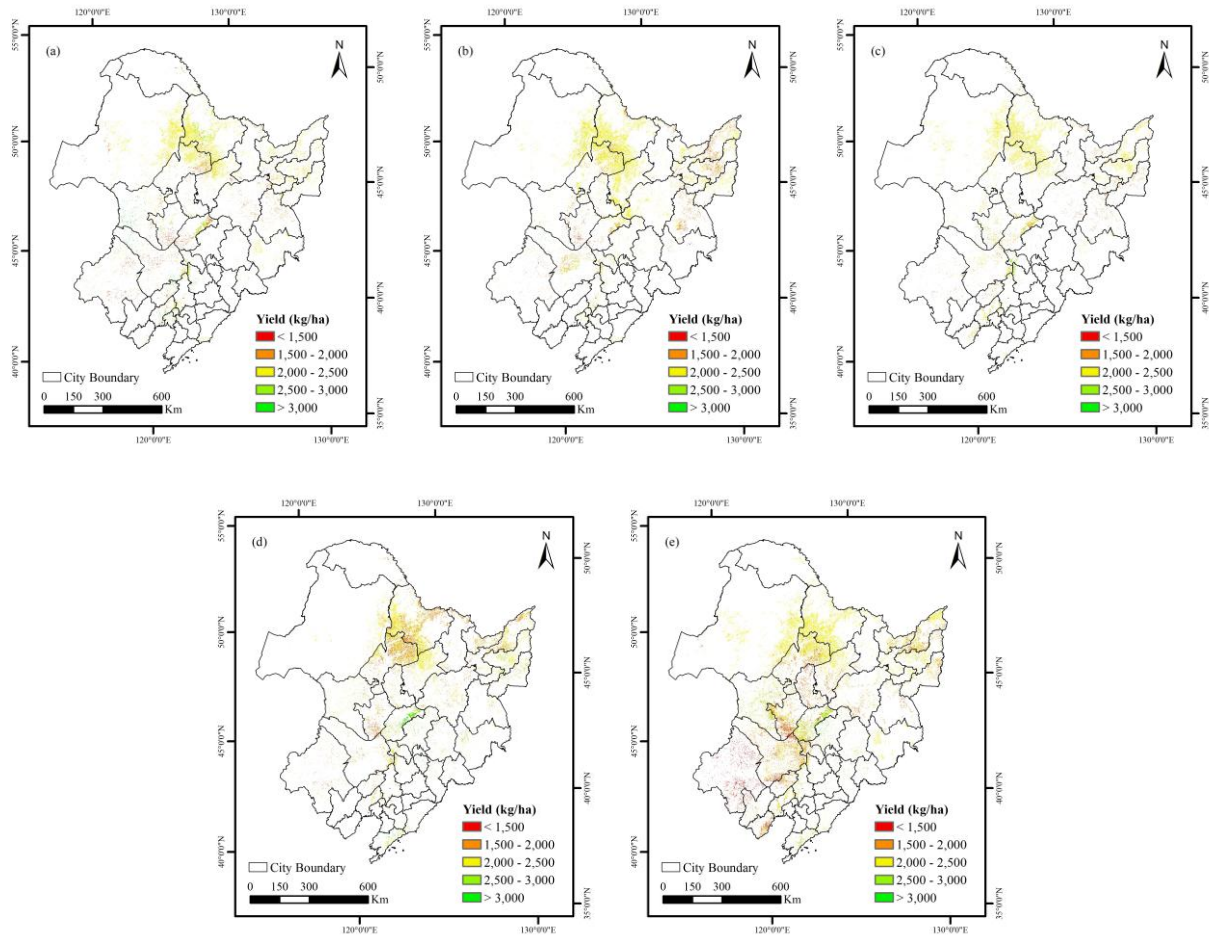
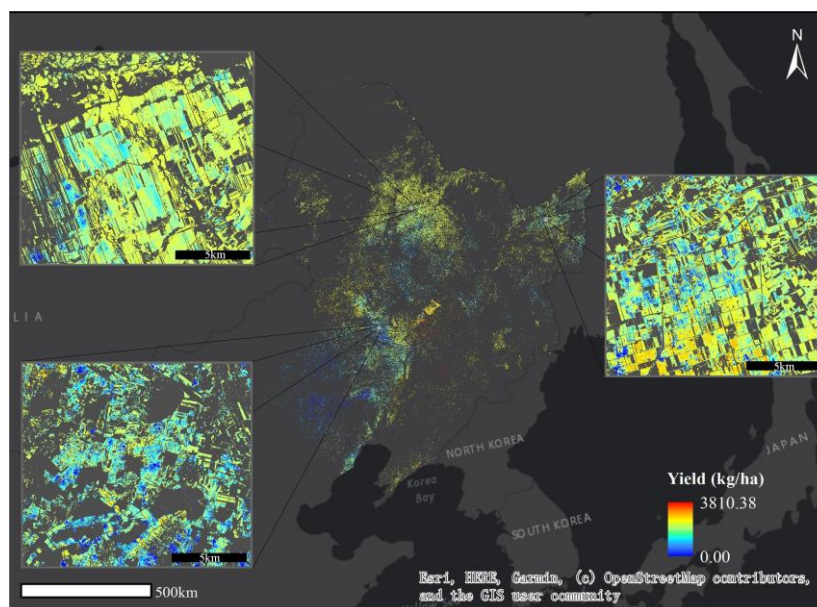


Figure 8: Spatial distribution of annual soybean yield derived from Sentinel-2 after calibration in Northeast China from 2019 to 2023.



390

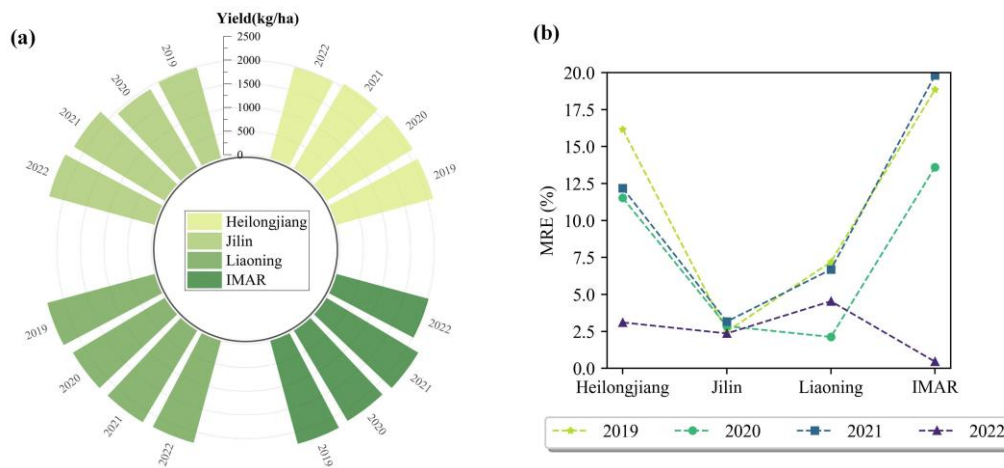
Figure 9: An example of yield estimation result (for the year 2023) used to showcase detailed local estimates.

Table 4 Mean values of statistical soybean yield at municipal scale in Northeast China compared with mean values, standard deviation (STD) and coefficient of variation (CV) for estimated soybean yield in Northeast China.

Year	Statistics (kg ha ⁻¹)	Mean (kg ha ⁻¹)	STD (kg ha ⁻¹)	CV (%)
2019	2137.24	2150.02	504.61	23.47
2020	2069.08	2125.49	372.21	17.51
2021	2115.57	2136.65	374.58	17.53
2022	2073.68	2036.89	465.26	22.84
2023	—	2035.34	603.43	29.65

We further analysed the spatial-temporal variation of soybean yield at the provincial scale (Fig. 10). On average, provincial-scale estimates achieved a mean relative error (MRE) of 7.94 % (Fig. 10b), with the highest accuracy observed in 2022 (Fig. 10b), mirroring municipal-level results (Fig. 6d). Over 2019 – 2022, Liaoning Province consistently exhibited the highest yields, whereas Heilongjiang Province, despite having the largest soybean planting area, recorded the lowest yields (Fig. 10a). This disparity likely due to Heilongjiang’s cooler climate, where later planting date result in shorter soybean development length. Across the four provinces, yields remained relatively stable, except in Jilin Province showed greatest interannual fluctuation exhibiting distinct decline followed by recovery. These results underscore the capacity of the proposed hybrid framework to capture spatial-temporal variations in soybean production.

400

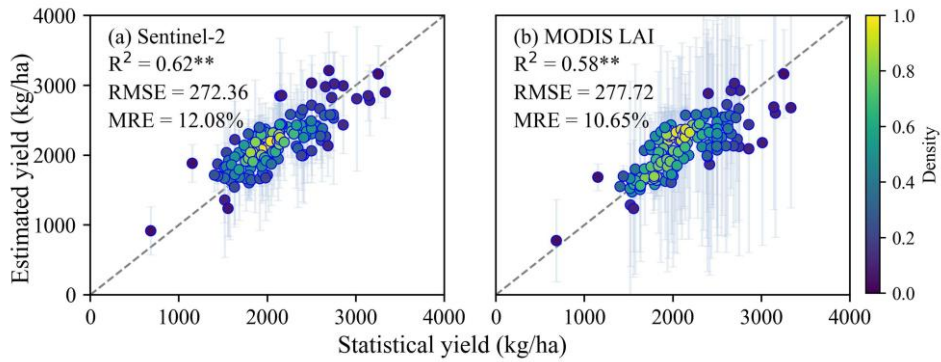


405 **Figure 10: Accuracy of the soybean yield estimation at provincial scale in Northeast China from 2019 to 2022. (a) represents the change in estimated yield for each province through the years; (b) represents MRE of results compared with statistical yield for each province**

5 Discussion

5.1 The complementarity between MODIS and Sentinel-2

410 This study generated soybean yield estimates using both MODIS LAI (500 m) products and S2 derived LAI (20 m) data. Over 2019 – 2022, the MODIS-based estimates achieved an overall R^2 of 0.58 ($p < 0.01$), an RMSE of $272.36 \text{ kg ha}^{-1}$ and an MRE of 12.08 % (Fig. 11b), slightly lower than the Sentinel-2 based results (Fig. 11a). The uncertainty of MODIS based estimates was higher than that the Sentinel-2 based estimates, likely reflecting MODIS’s coarser resolution. However, the Sentinel-2 based estimates exhibit inherent seaming effects caused by cloud-affected tile edges. We additionally used MODIS LAI to bias-correct Sentinel-2 yield maps, effectively minimizing the striping (“seaming”) effects in the 20 m products (Fig. 9), while preserving pixel-level detail through tile-based calibration (Fig. 13). Despite difference in spatial resolution, both MODIS and 415 Sentinel-2 satellite data demonstrated comparable ability to capture spatiotemporal variation in soybean yield (Fig. 12), achieving correlations with statistical data > 0.55 and overall errors $< 13 \%$ across all years.



420 **Figure 11: Comparison between estimated and statistical yield for 2019 – 2022 using Sentinel-2 (a) and MODIS LAI (b), respectively (excluding 2023 for which no statistical data was reported). The error-bars represent one standard deviation indicating the uncertainty of yield estimations. Dashed line represents 1:1 line. ** denotes statistical significance at $p < 0.01$.**

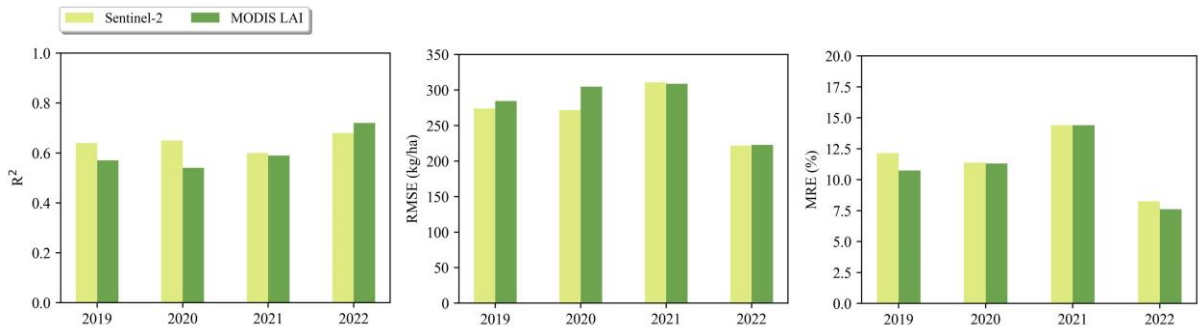


Figure 12: Comparison of accuracy evaluation results for soybean yield estimation in 2019 – 2022 (excluding 2023, for which no statistical data was reported) using Sentinel-2 and MODIS LAI data, respectively.

425 In practical applications, balancing both temporal and spatial resolution is critical for achieving robust yield prediction results (Azzari et al., 2017). Figure 13 compares the Sentinel-2 yield maps and the MODIS LAI yield maps within a 10 km grid under different soybean coverage. Thanks to 4-day revisit, MODIS LAI provides more cloud-free observations during the critical growth stages, improving the reliability of two LAI metrics (LAI_{mean1} and LAI_{mean2}). Its coarser spatial resolution also accelerates spatial processing over large areas. However, Sentinel-2's finer more effectively resolves intra-field yield
 430 heterogeneity (Fig. 13). MODIS-derived maps occasionally underestimated yields due to mixed pixels containing non-crop features (e.g., infrastructure), whereas Sentinel-2 minimized such errors.

While this study prioritized high-resolution mapping (using MODIS solely for Sentinel-2 seam correction), combining high spatial data (e.g., Sentinel-2 or UAV imagery) with high temporal frequency satellites (e.g., geostationary sensors or radar) could provide an optimal data source for crop-yield modelling (Gao and Anderson, 2019; He et al., 2018).

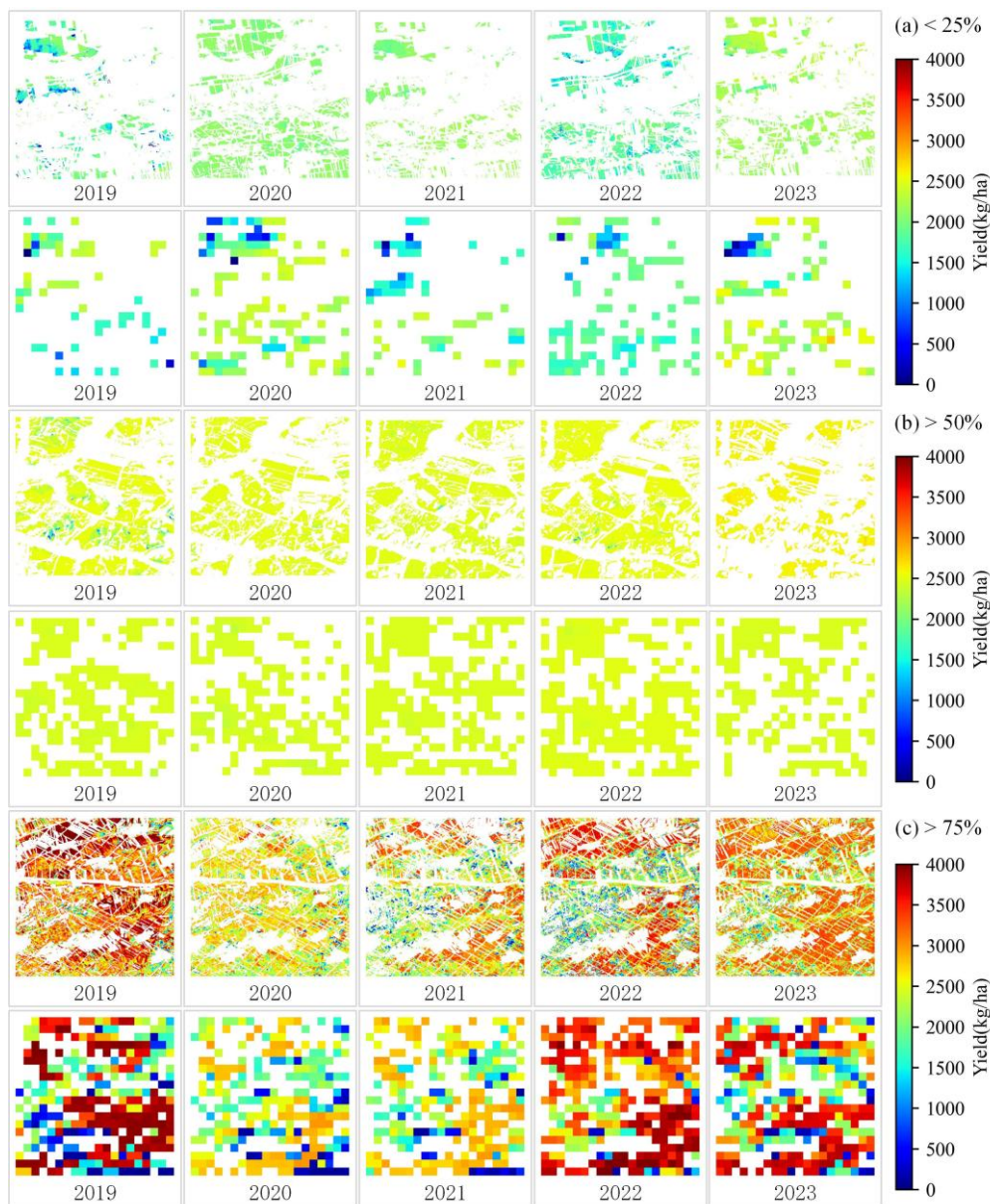


Figure 13: Comparisons of soybean yield estimation within a 10 km grid under different soybean coverage using Sentinel-2 (20 m) and MODIS LAI (500 m) data, where (a), (b), (c) represent soybean coverage less than 25%, more than 50% and more than 75%, respectively.

5.2 Advancements in this study

440 Accurate monitoring of soybean yield is crucial for food policy decision-making and security assessment. While previous studies have primarily explored the impact of environmental factors such as climate on soybean productivity (Guo et al., 2022; Zhao et al., 2023a), few efforts have focused on producing high-resolution soybean yield dataset for China's major soybean-

producing regions. To address this gap, our study produced the NortheastChinaSoybeanYield20m dataset, a 20-meter resolution dataset generated through a hybrid framework integrating the mechanistic WOFOST crop growth model and a GRU
445 deep learning algorithm. Unlike purely data-driven approaches that rely on extensive ground data, our approach leveraged both data mining capabilities and mechanistic modelling, which improve the model's interpretability and enhances its potential for transferability across regions. The integration of the WOFOST model ensured the simulation of diverse production scenarios under varying climate, soil, crop variety and management conditions, providing a robust synthetic training data for the GRU network. This combination allowed the model to generate well, even in areas with limited observational data, therefore
450 overcoming common limitations related to data scarcity and high computational costs. Accuracy assessments using both in-situ and statistical yield data confirmed that the generated NortheastChinaSoybeanYield20m dataset delivered reliable yield estimates across field and regional scales (Fig. 5 and 6). The results also verified the model's stability across time and space, reinforcing its potential for large-scale agricultural monitoring and strategic planning.

When compared to previous studies using integrated remote sensing data and process-based model to estimate soybean
455 yield, for instance, Baup et al., (2015) reported estimation error ranging from 2% to 18%, our method achieved comparable levels of accuracy. It also outperformed existing field-scale studies (e.g., RMSE = 400.946 kg ha⁻¹ in Ren et al., (2023) and MRE of 29.73% in Du et al., (2014)) and municipal-scale models (e.g., RMSE = 16 % in Von Bloh et al., (2023)). Furthermore, the NortheastChinaSoybeanYield20m dataset showed improved performance relative to similar high-resolution soybean yield products from other countries (e.g., annual 30 m soybean yield mapping in Brazil, with R² values between 0.31 and 0.71 and
460 RMSEs ranging from 275 to 740 kg ha⁻¹ (Song et al., 2022)).

Although studies based on UAV and RGB data have demonstrated even higher soybean yield estimation accuracy (Li et al., 2021, 2024), such methods are often constrained by high costs and limited spatial coverage, making them impractical for large-scale applications. In contrast, the method developed in this study offers a well-balanced solution that combines computational efficiency, high spatial resolution, and strong predictive accuracy. Our approach offers scalable and practical
465 solution for producing high-resolution, large-scale crop yield datasets.

5.3 Limitations and future developments

In this study, a multi-scenario soybean growth dataset was developed by simulating various combinations input parameters within the WOFOST model. These diverse scenarios were designed to reflect different environmental and management conditions, ultimately serving as training data for the yield estimation model. One advantage of the model is its scalability, it
470 can be readily applied to other regions and countries that lack sufficient ground observation data, such as parts of Africa and India, thus offering a promising tool for global agricultural monitoring.

However, the validation results revealed some notable limitations. Specifically, the model exhibited a tendency to produce large uncertainty in low- or high- yielding areas, introducing error into the overall yield estimation (Fig. 5 and 6). This pattern suggests a systematic bias in the model's predictions, particularly in regions with extreme yield values. Additionally, spatial
475 analysis showed that estimation errors were more pronounced in the northern region, where is characterized by complex terrain,

compared to the relatively flat central region (Fig. 7). These discrepancies highlight the need to refine parameterization for extreme yield conditions and integrate higher-resolution environmental drivers (e.g., terrain, localized weather).

On the one hand, the estimation errors may be attributed to the inherent limitations of the WOFOST model. As a process-based model, WOFOST simplifies its calculations for simulating physiological processes, which can hinder its ability to fully replicate the complex realities of soybean in the field. Factors, such as pest infestations, diseases, and abiotic stresses are either oversimplified or excluded (Gaso et al., 2024). These omissions can lead to systematic simulation errors, particularly under stress conditions that significantly affect crop yield. Moreover, the parameterization of the WOFOST model in this study purely relied on values from literature and existing dataset rather than local optimization. As a result, local variability because of farming practices, soil properties, and environmental conditions may not have been adequately captured. This lacks local optimization likely result in higher estimation error, especially in complex landscapes with sparse ground observations. To address these issues, future works incorporating field-specific parameters or advanced data assimilation techniques could help reduce bias and improve model accuracy across heterogeneous landscapes. Given the spatial variability in soybean growth within the study area, constructing ecological zones based on factors like climate, elevation, and management practices might provide a more targeted model approach. For instance, Huang et al., (2023) defined the ecological zones through using Theissen polygons derived from meteorological station locations. This zoning strategy could enhance the representativeness of the training data and reduce yield estimation uncertainties.

On the other hand, the estimation errors may stem from the overfitting of the GRU model. The GRU was trained on the multi-scenarios simulated dataset, a large number of simulations that included all available combinations (e.g., all meteorological data), which introduced a significant amount of redundant information. The redundancy not only potentially reduce the dataset's representativeness, but also increase the computational burden during model training. As a result, the trained GRU model may have become overly turned to specific temporal patterns in certain years, limiting its ability to generalize to other time period or regions with different growth conditions. This overfitting effect might result in large yield estimation errors across different years and regions, particularly in areas where soybean yields deviated significantly from the norm. To address these issues, refining the structure and composition of the training dataset, and removing redundant information would enhance the diversity and quality of the training inputs. One potential approach to reduce redundancy is through spatiotemporal clustering of various environmental (e.g., meteorological station data), which could filter out stations with highly similar information. Moreover, monitoring the validation error throughout the training process, and implementing regularization techniques (e.g., L2 weight regularization) could help to prevent overfitting and improve the GRU model's generalization capability, leading to improve soybean estimation across varying conditions.

Finally, accurate estimation of soybean yield depended on the quality of the input data. The spatial resolution of remote sensing imagery could limit the model's ability to predict spatial variability in yield. In this study, ERA5-land dataset was applied to obtain the spatial-temporal distribution of soybean phenology in the study area. To be consistent with Sentinel-2 data, all datasets were resampled to a 20 m resolution. Downscaling the coarse spatial data could increase the uncertainty of inputs to the model.

510 Moreover, the issue of mixed pixels led to a loss of detailed surface information (Zhao et al., 2023), particularly in
heterogeneous or complex environments. With the advent of higher spatial-temporal resolution remote sensing data, the
estimation accuracy of crop yield is expected to be further improved. The mean values of LAI at two key soybean growth
stages were used as the primary remote sensing-based input features for yield estimation. However, errors in LAI retrieval
515 from remote sensing data also contributed uncertainty in yield predictions. Integrating agronomic knowledge with remote
sensing mechanisms has emerged as a promising way to reduce the uncertainty and improve model reliability (Chen et al.,
2022; Hu et al., 2024). Coupling radiative transfer model such as PROSAIL (Jacquemoud et al., 2009)) with crop growth
model can enhance the simulation of leaf and canopy characteristics and provide additional constraints for more accurate
modelling (Ntakos et al., 2024). In addition, the combination of IoT, blockchain, and precision agriculture with machine
520 learning and biophysical models can offer a powerful framework for sustainable agricultural monitoring, addressing challenges
in data heterogeneity, model scalability, and decision-making processes. These technologies can facilitate real-time data
collection, ensure data security and transparency. Precision agriculture techniques, combined with advanced sensing
technologies, can effectively improve the accuracy and timeliness of input data, addressing current limitations in model
calibration, validation and prediction.

6 Data availability

525 The soybean yield dataset for Northeast China (NortheastChinaSoybeanYield20m) during the 2019-2023 period is available
at <https://doi.org/10.5281/zenodo.14263103> (Xu et al., 2024).

7 Conclusions

This study generated a high-resolution (20 m) soybean yield dataset for Northeast China from 2019 to 2023
(NortheastChinaSoybeanYield20m) using a hybrid framework that couple the WOFOST crop growth model with a Gated
530 Recurrent Unit (GRU) deep learning algorithm. The framework leveraged a comprehensive soybean growth dataset simulated
by WOFOST, which accounted for diverse production scenarios, including variations in climates, crop varieties, soil types and
agro-managements practices. This approach effectively reduces reliance on ground observation data, which demonstrating
enhanced spatiotemporal generalization capabilities.

The dataset was conducted using multi-source remote sensing data, with Sentinel-2 derived time-series LAI as the primary
535 input. Yield estimations showed robust performance at both field and municipal scales, achieving RMSE of 287.44 kg ha⁻¹ and
272.36 kg ha⁻¹, respectively. To address spatial discontinuities in Sentinel-2 data, corrections using MODIS LAI-derived yield
maps effectively mitigated seam effects, achieving complementary benefits in temporal and spatial resolution. The final dataset
exhibits high temporal stability and spatial continuity, with mean relative errors (MRE) averaging of 11.46 % at the municipal
scale and 7.94 % at the provincial scale.

540 The NortheastChinaSoybeanYield20m dataset successfully captures fine-scale spatiotemporal variations in soybean yield, offering potentials for optimizing production strategies, guiding precision agriculture, and enhancing food security and policy.

Authorship contributions

JX (first author) and QL – conceptualization; JX (first author), XD, YZ, HW, JX, YS and YD – data curation; JX (first author), XD, TD – methodology; JX (first author), XD, JX and JZ – investigation; TD and QL – supervision; HW, JX (first author) and JZ – validation; YZ, HW and JZ – visualization; JX (first author) – original draft preparation; XD, TD and YZ – reviewing and editing the manuscript.

Competing interests

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

550 Acknowledgements

This research was funded by the National Key R&D Program of China (2021YFD1500103), the Strategic Priority Research Program of the Chinese Academy of Sciences (XDA28070504), the National Science Foundation of China (42371359), and the Key Program of High-resolution Earth Observation System (71-Y50G10-9001-22/23).

Appendix A

555 Table A1 Values of crop parameters in WOFOST.

Parameter	Description	Units	Value	Source
Crop initial parameters				
TDWI	Initial total crop dry weight	kg ha ⁻¹	120	Default value in WOFOST
RGRLAI	Maximum relative increase in LAI	ha ha ⁻¹ d ⁻¹	0.01	Default value in WOFOST
Parameters for emergence				
TBASEM	Minimum threshold temperature for emergence	°C	8.0	Qu et al., (2023)
TEFFMX	Maximum threshold temperature for emergence	°C	22.0	Qu et al., (2023)

TSUMEM	Accumulated temperature from sowing to emergence	°C	70.0	Qu et al., (2023)
Phenological parameters				
DLO	Optimal daylength for development	h	-99	Default value in WOFOST
DLC	Critical daylength	h	-99	Default value in WOFOST
TSUM1	Cumulative temperature from emergence to anthesis	°C	450 (early maturity) 480 (medium-early maturity) 520 (intermediate maturity) 540 (medium-late maturity) 580 (late maturity)	Qu et al., (2023)
TSUM2	Cumulative temperature from anthesis to maturity	°C	660 (early maturity) 770 (medium-early maturity) 870 (intermediate maturity) 960 (medium-late maturity) 1000 (late maturity)	Qu et al., (2023)
Green area parameters				
TBASE	Lower threshold temperature for aging of leaves	°C	7.0	Default value in WOFOST
SPAN	Life span of leaves growing at 35 °C	d	23	Default value in WOFOST
SLATB00	Specific leaf area at DVS = 0.00	ha kg ⁻¹	0.00140	Default value in WOFOST
SLATB045	Specific leaf area at DVS = 0.45	ha kg ⁻¹	0.00250	Default value in WOFOST
SLATB090	Specific leaf area at DVS = 0.90	ha kg ⁻¹	0.00250	Default value in WOFOST
SLATB200	Specific leaf area at DVS = 2.00	ha kg ⁻¹	0.00070	Default value in WOFOST
Assimilation parameters				
KDIFTB00	Extinction coefficient for diffuse visible light (DVS = 0)	-	0.80	Default value in WOFOST

KDIFTB200	Extinction coefficient for diffuse visible light (DVS = 2)	-	0.80	Default value in WOFOST
EFFTB0	Light use efficiency of a single leaf (T = 0 °C)	kg ha ⁻¹ h ⁻¹ J ⁻¹ m ² s ⁻¹	0.40	Default value in WOFOST
EFTB40	Light use efficiency of a single leaf (T = 40 °C)	kg ha ⁻¹ h ⁻¹ J ⁻¹ m ² s ⁻¹	0.40	Default value in WOFOST
AMAXTB00	Maximum leaf CO ₂ assimilation rate (DVS = 0)	kg ha ⁻¹ h ⁻¹	29.00	Default value in WOFOST
AMAXTB170	Maximum leaf CO ₂ assimilation rate (DVS = 1.7)	kg ha ⁻¹ h ⁻¹	25.31	Sun et al., (2022)
AMAXTB200	Maximum leaf CO ₂ assimilation rate (DVS = 2)	kg ha ⁻¹ h ⁻¹	0.00	Default value in WOFOST
TMPFTB00	Reduction factor of AMAX (T = 0 °C)	-	0.00	Default value in WOFOST
TMPFTB10	Reduction factor of AMAX (T = 10 °C)	-	0.30	Default value in WOFOST
TMPFTB20	Reduction factor of AMAX (T = 20 °C)	-	0.60	Default value in WOFOST
TMPFTB25	Reduction factor of AMAX (T = 25 °C)	-	0.80	Default value in WOFOST
TMPFTB30	Reduction factor of AMAX (T = 30 °C)	-	1.00	Default value in WOFOST
TMPFTB35	Reduction factor of AMAX (T = 35 °C)	-	1.00	Default value in WOFOST
Conversion of assimilates into biomass				
CVL	Conversion efficiency of assimilates into leaf tissue	kg kg ⁻¹	0.72	Default value in WOFOST
CVO	Conversion efficiency of assimilates into storage organs	kg kg ⁻¹	0.48	Default value in WOFOST
CVR	Conversion efficiency of assimilates into root tissue	kg kg ⁻¹	0.72	Default value in WOFOST
CVS	Conversion efficiency of assimilates into stem tissue	kg kg ⁻¹	0.69	Default value in WOFOST

Maintenance respiration parameters

Q10	Relative change in respiration rate per 10 °C temperature increase	-	2.0	Default value in WOFOST
RML	Relative maintenance respiration rate of leaves	kg CH ₂ O kg ⁻¹ d ⁻¹	0.03	Default value in WOFOST
RMO	Relative maintenance respiration rate of storage organs	kg CH ₂ O kg ⁻¹ d ⁻¹	0.017	Default value in WOFOST
RMR	Relative maintenance respiration rate of roots	kg CH ₂ O kg ⁻¹ d ⁻¹	0.01	Default value in WOFOST
RMS	Relative maintenance respiration rate of stems	kg CH ₂ O kg ⁻¹ d ⁻¹	0.015	Default value in WOFOST

Partitioning parameters

FRTB00	Fraction of total dry matter to roots at DVS = 0	kg kg ⁻¹	0.62	Sun et al., (2022)
FRTB075	Fraction of total dry matter to roots at DVS = 0.75	kg kg ⁻¹	0.35	Default value in WOFOST
FRTB100	Fraction of total dry matter to roots at DVS = 1	kg kg ⁻¹	0.15	Default value in WOFOST
FRTB150	Fraction of total dry matter to roots at DVS = 1.5	kg kg ⁻¹	0.00	Default value in WOFOST
FRTB200	Fraction of total dry matter to roots at DVS = 2.0	kg kg ⁻¹	0.00	Default value in WOFOST
FLTB00	Fraction of total dry matter to leaves at DVS = 0	kg kg ⁻¹	0.70	Default value in WOFOST
FLTB100	Fraction of total dry matter to leaves at DVS = 1.0	kg kg ⁻¹	0.70	Default value in WOFOST
FLTB115	Fraction of total dry matter to leaves at DVS = 1.15	kg kg ⁻¹	0.60	Default value in WOFOST
FLTB130	Fraction of total dry matter to leaves at DVS = 1.3	kg kg ⁻¹	0.43	Default value in WOFOST
FLTB150	Fraction of total dry matter to leaves at DVS = 1.5	kg kg ⁻¹	0.15	Default value in WOFOST

FLTB200	Fraction of total dry matter to leaves at DVS = 2.0	kg kg ⁻¹	0.00	Default value in WOFOST
FSTB00	Fraction of total dry matter to stems at DVS = 0	kg kg ⁻¹	0.30	Default value in WOFOST
FSTB100	Fraction of total dry matter to stems at DVS = 1.0	kg kg ⁻¹	0.30	Default value in WOFOST
FSTB115	Fraction of total dry matter to stems at DVS = 1.15	kg kg ⁻¹	0.25	Default value in WOFOST
FSTB130	Fraction of total dry matter to stems at DVS = 1.3	kg kg ⁻¹	0.10	Default value in WOFOST
FSTB150	Fraction of total dry matter to stems at DVS = 1.5	kg kg ⁻¹	0.10	Default value in WOFOST
FSTB200	Fraction of total dry matter to stems at DVS = 2.0	kg kg ⁻¹	0.00	Default value in WOFOST
FOTB00	Fraction of total dry matter to storage organs at DVS = 0	kg kg ⁻¹	0.00	Default value in WOFOST
FOTB100	Fraction of total dry matter to storage organs at DVS = 1.0	kg kg ⁻¹	0.00	Default value in WOFOST
FOTB115	Fraction of total dry matter to storage organs at DVS = 1.15	kg kg ⁻¹	0.15	Default value in WOFOST
FOTB130	Fraction of total dry matter to storage organs at DVS = 1.3	kg kg ⁻¹	0.47	Default value in WOFOST
FOTB150	Fraction of total dry matter to storage organs at DVS = 1.5	kg kg ⁻¹	0.75	Default value in WOFOST
FOTB200	Fraction of total dry matter to storage organs at DVS = 2.0	kg kg ⁻¹	1.00	Default value in WOFOST
Death rate parameters				
PERDL	Maximum relative death rate of leaves due to water stress	kg kg ⁻¹ d ⁻¹	0.03	Default value in WOFOST
RDRRTB00	Relative death rate of roots at DVS = 0	kg kg ⁻¹ d ⁻¹	0.00	Default value in WOFOST
RDRRTB150	Relative death rate of roots at DVS = 1.5	kg kg ⁻¹ d ⁻¹	0.00	Default value in WOFOST

RDRRTB151	Relative death rate of roots at DVS = 1.51	kg kg ⁻¹ d ⁻¹	0.02	Default value in WOFOST
RDRRTB200	Relative death rate of roots at DVS = 2.0	kg kg ⁻¹ d ⁻¹	0.02	Default value in WOFOST
RDRSTB00	Relative death rate of stems at DVS = 0	kg kg ⁻¹ d ⁻¹	0.00	Default value in WOFOST
RDRSTB150	Relative death rate of stems at DVS = 1.5	kg kg ⁻¹ d ⁻¹	0.00	Default value in WOFOST
RDRSTB151	Relative death rate of stems at DVS = 1.51	kg kg ⁻¹ d ⁻¹	0.02	Default value in WOFOST
RDRSTB200	Relative death rate of stems at DVS = 2.0	kg kg ⁻¹ d ⁻¹	0.02	Default value in WOFOST
Water use parameters				
CFET	Correction factor transpiration rate	-	1.0	Default value in WOFOST
DEPNR	Crop group number for soil water depletion	-	5.0	Default value in WOFOST
IAIRDU	Air ducts in roots present (=1) or not (=0)	-	0	Default value in WOFOST
IOX	Oxygen stress effect enabled (=1) or not (=0)	-	0	Default value in WOFOST
Rooting parameters				
RDI	Initial rooting depth	cm	10	Default value in WOFOST
RRI	Maximum daily increase in rooting depth	cm d ⁻¹	1.2	Default value in WOFOST
RDMCR	Maximum rooting depth	cm	120	Default value in WOFOST

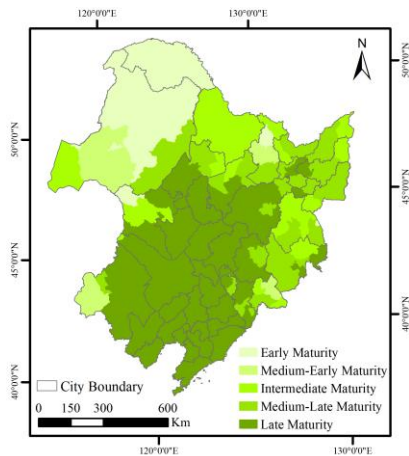


Figure A1: Spatial distribution of soybean types in Northeast China.

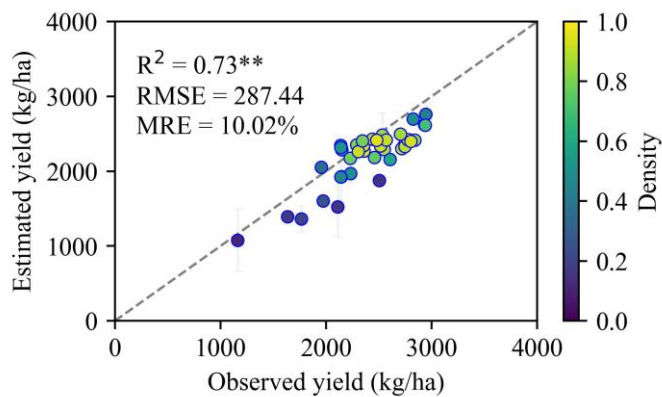
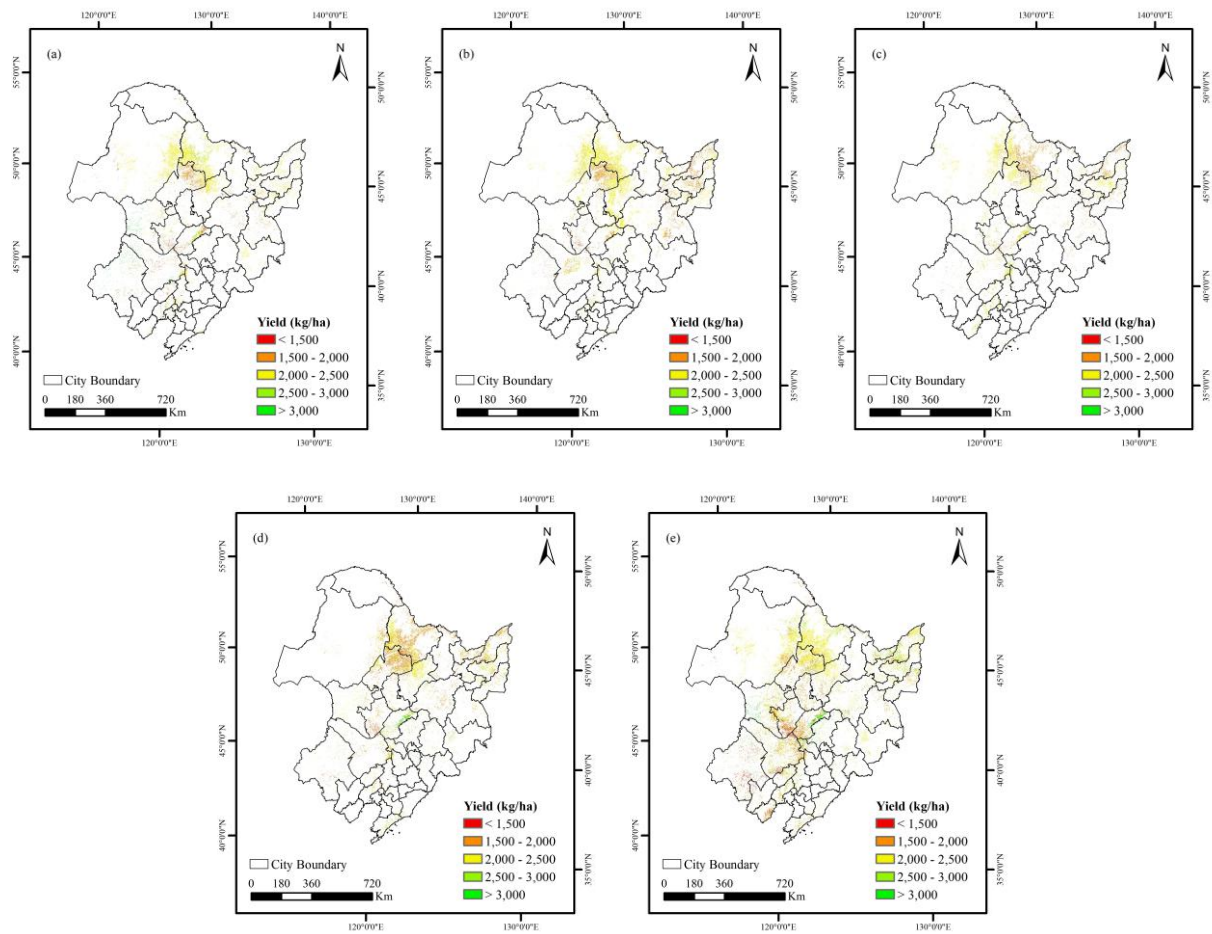


Figure A2: Comparison between estimated and observed yield for both of 2022 and 2023. The error-bars represent one standard deviation indicating the uncertainty of yield estimations. Dashed lines represent 1:1 line. ** denotes statistical significance at $p < 0.01$.



565 **Figure A3: Spatial distribution of annual soybean yield derived from Sentinel-2 before calibration in Northeast China from 2019 to 2023.**

References

- Açikkar, M.: Fast grid search: A grid search-inspired algorithm for optimizing hyperparameters of support vector regression, *Turkish Journal of Electrical Engineering and Computer Sciences*, 32, 68–92, <https://doi.org/10.55730/1300-0632.4056>, 2024.
- 570 Allen, L. H., Kirkham, M. B., Olszyk, D. M., Whitman, C. E., and Pickering, N. B.: Plant Modeling: Advances and Gaps in Our Capability to Predict Future Crop Growth and Yield in Response to global Climate Change, *Advances in Carbon Dioxide Effects Research*, 1997.
- Ang, Y., Shafri, H. Z. M., Lee, Y. P., Abidin, H., Bakar, S. A., Hashim, S. J., Che'Ya, N. N., Hassan, M. R., Lim, H. S., and Abdullah, R.: A novel ensemble machine learning and time series approach for oil palm yield prediction using Landsat time series imagery based on NDVI, *Geocarto International*, 37, 9865–9896, <https://doi.org/10.1080/10106049.2022.2025920>, 2022.
- 575

- Azzari, G., Jain, M., and Lobell, D. B.: Towards fine resolution global maps of crop yields: Testing multiple methods and satellites in three countries, *Remote Sensing of Environment*, 202, 129–141, <https://doi.org/10.1016/j.rse.2017.04.014>, 2017.
- 580 Baup, F., Fieuzal, R., and Betbeder, J.: Estimation of soybean yield from assimilated optical and radar data into a simplified agrometeorological model, in: 2015 IEEE International Geoscience and Remote Sensing Symposium (IGARSS), IGARSS 2015 - 2015 IEEE International Geoscience and Remote Sensing Symposium, Milan, Italy, 3961–3964, <https://doi.org/10.1109/IGARSS.2015.7326692>, 2015.
- Cao, J., Zhang, Z., Tao, F., Zhang, L., Luo, Y., Zhang, J., Han, J., and Xie, J.: Integrating Multi-Source Data for Rice Yield Prediction across China using Machine Learning and Deep Learning Approaches, *Agricultural and Forest Meteorology*, 297, 108275, <https://doi.org/10.1016/j.agrformet.2020.108275>, 2021.
- 585 Chen, Q., Zheng, B., Chen, T., and Chapman, S. C.: Integrating a crop growth model and radiative transfer model to improve estimation of crop traits based on deep learning, *Journal of Experimental Botany*, 73, 6558–6574, <https://doi.org/10.1093/jxb/erac291>, 2022.
- Chen, Y., Liu, S., Li, H., Li, X. F., Song, C. Y., Cruse, R. M., and Zhang, X. Y.: Effects of conservation tillage on corn and soybean yield in the humid continental climate region of Northeast China, *Soil and Tillage Research*, 115–116, 56–61, <https://doi.org/10.1016/j.still.2011.06.007>, 2011.
- 590 Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y.: Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation, in: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), Doha, Qatar, 1724–1734, <https://doi.org/10.3115/v1/D14-1179>, 2014.
- 595 Choi, D.-H., Ban, H.-Y., Seo, B.-S., Lee, K.-J., and Lee, B.-W.: Phenology and Seed Yield Performance of Determinate Soybean Cultivars Grown at Elevated Temperatures in a Temperate Region, *PLoS ONE*, 11, e0165977, <https://doi.org/10.1371/journal.pone.0165977>, 2016.
- Diepen, C. A., Wolf, J., Keulen, H., and Rappoldt, C.: WOFOST: a simulation model of crop production, *Soil Use & Management*, 5, 16–24, <https://doi.org/10.1111/j.1475-2743.1989.tb00755.x>, 1989.
- 600 Dokoohaki, H., Kivi, M. S., Martinez-Feria, R., Miguez, F. E., and Hoogenboom, G.: A comprehensive uncertainty quantification of large-scale process-based crop modeling frameworks, *Environ. Res. Lett.*, 16, 084010, <https://doi.org/10.1088/1748-9326/ac0f26>, 2021.
- Dong, T., Liu, J., Qian, B., He, L., Liu, J., Wang, R., Jing, Q., Champagne, C., McNairn, H., Powers, J., Shi, Y., Chen, J. M., and Shang, J.: Estimating crop biomass using leaf area index derived from Landsat 8 and Sentinel-2 data, *ISPRS Journal of Photogrammetry and Remote Sensing*, 168, 236–250, <https://doi.org/10.1016/j.isprsjprs.2020.08.003>, 2020.
- 605 Du, X., Song, F., Wang, H., Huanxuezhang, Meng, J., Li, Q., Liu, J., Ding, L., and Lu, Y.: Soybean yield estimation using HJ-1 CCD data in Northeast China, in: 2014 The Third International Conference on Agro-Geoinformatics, 2014 Third International Conference on Agro-Geoinformatics, Beijing, China, 1–4, <https://doi.org/10.1109/Agro-Geoinformatics.2014.6910627>, 2014.
- 610 Du, X., Zhu, J., Xu, J., Li, Q., Tao, Z., Zhang, Y., Wang, H., and Hu, H.: Remote sensing-based winter wheat yield estimation integrating machine learning and crop growth multi-scenario simulations, *International Journal of Digital Earth*, 18, 2443470, <https://doi.org/10.1080/17538947.2024.2443470>, 2025.

- 615 Duchemin, B., Maisongrande, P., Boulet, G., and Benhadj, I.: A simple algorithm for yield estimates: Evaluation for semi-arid irrigated winter wheat monitored with green leaf area index, *Environmental Modelling & Software*, 23, 876–892, <https://doi.org/10.1016/j.envsoft.2007.10.003>, 2008.
- Falcon, W. P., Naylor, R. L., and Shankar, N. D.: Rethinking Global Food Demand for 2050, *Population & Development Rev*, 48, 921–957, <https://doi.org/10.1111/padr.12508>, 2022.
- Fan, R., Zhang, X., Liang, A., Shi, X., Chen, X., Bao, K., Yang, X., and Jia, S.: Tillage and rotation effects on crop yield and profitability on a Black soil in northeast China, *Can. J. Soil. Sci.*, 92, 463–470, <https://doi.org/10.4141/cjss2010-020>, 2012.
- 620 FAOSTAT: FAO statistical database, 2022.
- Feng, P., Wang, B., Liu, D. L., Waters, C., Xiao, D., Shi, L., and Yu, Q.: Dynamic wheat yield forecasts are improved by a hybrid approach using a biophysical model and machine learning technique, *Agricultural and Forest Meteorology*, 285–286, 107922, <https://doi.org/10.1016/j.agrformet.2020.107922>, 2020.
- 625 Gao, F. and Anderson, M.: Evaluating Yield Variability of Corn and Soybean Using Landsat-8, Sentinel-2 and Modis in Google Earth Engine, in: *IGARSS 2019 - 2019 IEEE International Geoscience and Remote Sensing Symposium, IGARSS 2019 - 2019 IEEE International Geoscience and Remote Sensing Symposium*, Yokohama, Japan, 7286–7289, <https://doi.org/10.1109/IGARSS.2019.8897990>, 2019.
- 630 Gaso, D. V., Paudel, D., De Wit, A., Puntel, L. A., Mullissa, A., and Kooistra, L.: Beyond assimilation of leaf area index: Leveraging additional spectral information using machine learning for site-specific soybean yield prediction, *Agricultural and Forest Meteorology*, 351, 110022, <https://doi.org/10.1016/j.agrformet.2024.110022>, 2024.
- Gevaert, C. M.: Explainable AI for earth observation: A review including societal and regulatory perspectives, *International Journal of Applied Earth Observation and Geoinformation*, 112, 102869, <https://doi.org/10.1016/j.jag.2022.102869>, 2022.
- 635 Gitelson, A. and Merzlyak, M. N.: Spectral Reflectance Changes Associated with Autumn Senescence of *Aesculus hippocastanum* L. and *Acer platanoides* L. Leaves. Spectral Features and Relation to Chlorophyll Estimation, *Journal of Plant Physiology*, 143, 286–292, [https://doi.org/10.1016/S0176-1617\(11\)81633-0](https://doi.org/10.1016/S0176-1617(11)81633-0), 1994.
- Gopi, P. S. S. and Karthikeyan, M.: Red fox optimization with ensemble recurrent neural network for crop recommendation and yield prediction model, *Multimed Tools Appl*, 83, 13159–13179, <https://doi.org/10.1007/s11042-023-16113-2>, 2023.
- Graham, P. H. and Vance, C. P.: Legumes: Importance and constraints to greater use, *Plant physiology*, 131, 872–877, 2003.
- 640 Guo, S., Guo, E., Zhang, Z., Dong, M., Wang, X., Fu, Z., Guan, K., Zhang, W., Zhang, W., Zhao, J., Liu, Z., Zhao, C., and Yang, X.: Impacts of mean climate and extreme climate indices on soybean yield and yield components in Northeast China, *Science of The Total Environment*, 838, 156284, <https://doi.org/10.1016/j.scitotenv.2022.156284>, 2022.
- He, M., Kimball, J., Maneta, M., Maxwell, B., Moreno, A., Beguería, S., and Wu, X.: Regional Crop Gross Primary Productivity and Yield Estimation Using Fused Landsat-MODIS Data, *Remote Sensing*, 10, 372, <https://doi.org/10.3390/rs10030372>, 2018.
- 645 Hu, P., Zheng, B., Chen, Q., Grunefeld, S., Choudhury, M. R., Fernandez, J., Potgieter, A., and Chapman, S. C.: Estimating aboveground biomass dynamics of wheat at small spatial scale by integrating crop growth and radiative transfer models with satellite remote sensing data, *Remote Sensing of Environment*, 311, 114277, <https://doi.org/10.1016/j.rse.2024.114277>, 2024.

- Huang, H., Huang, J., Wu, Y., Zhuo, W., Song, J., Li, X., Li, L., Su, W., Ma, H., and Liang, S.: The Improved Winter Wheat Yield Estimation by Assimilating GLASS LAI Into a Crop Growth Model With the Proposed Bayesian Posterior-Based Ensemble Kalman Filter, *IEEE Trans. Geosci. Remote Sensing*, 61, 1–18, <https://doi.org/10.1109/TGRS.2023.3259742>, 2023.
- 650 Huang, J., Tian, L., Liang, S., Ma, H., Becker-Reshef, I., Huang, Y., Su, W., Zhang, X., Zhu, D., and Wu, W.: Improving winter wheat yield estimation by assimilation of the leaf area index from Landsat TM and MODIS data into the WOFOST model, *Agricultural and Forest Meteorology*, 204, 106–121, <https://doi.org/10.1016/j.agrformet.2015.02.001>, 2015.
- Huang, J., Gómez-Dans, J. L., Huang, H., Ma, H., Wu, Q., Lewis, P. E., Liang, S., Chen, Z., Xue, J.-H., Wu, Y., Zhao, F., Wang, J., and Xie, X.: Assimilation of remote sensing into crop growth models: Current status and perspectives, *Agricultural and Forest Meteorology*, 276–277, 107609, <https://doi.org/10.1016/j.agrformet.2019.06.008>, 2019.
- 655 Huang, J., Song, J., Huang, H., Zhuo, W., Niu, Q., Wu, S., Ma, H., and Liang, S.: Progress and perspectives in data assimilation algorithms for remote sensing and crop growth model, *Science of Remote Sensing*, 10, 100146, <https://doi.org/10.1016/j.srs.2024.100146>, 2024.
- Huang, Y. and Liu, Z.: Improving Northeast China’s soybean and maize planting structure through subsidy optimization considering climate change and comparative economic benefit, *Land Use Policy*, 146, 107319, <https://doi.org/10.1016/j.landusepol.2024.107319>, 2024.
- 660 Hunt, M. L., Blackburn, G. A., Carrasco, L., Redhead, J. W., and Rowland, C. S.: High resolution wheat yield mapping using Sentinel-2, *Remote Sensing of Environment*, 233, 111410, <https://doi.org/10.1016/j.rse.2019.111410>, 2019.
- Jacquemoud, S., Verhoef, W., Baret, F., Bacour, C., Zarco-Tejada, P. J., Asner, G. P., François, C., and Ustin, S. L.: PROSPECT+SAIL models: A review of use for vegetation characterization, *Remote Sensing of Environment*, 113, S56–S66, <https://doi.org/10.1016/j.rse.2008.01.026>, 2009.
- 665 Jain, A. K. and Dubes, R. C.: Algorithms for clustering data, *Technometrics*, 32, 227–229, 1988.
- Kaur, S. and Singh, M.: Modeling the crop growth - A review, *MAUSAM*, 71, 103–114, 2020.
- Knyazikhin, Y. ;, Glassy, J. ;, Privette, J. L. ;, Tian, Y. ;, and Running, S. W. ; MODIS Leaf Area Index (LAI) and Fraction of Photosynthetically Active Radiation Absorbed by Vegetation (FPAR) Product (MOD15) Algorithm Theoretical Basis Document, 2018.
- 670 Li, C., Ma, C., Cui, Y., Lu, G., and Wei, F.: UAV Hyperspectral Remote Sensing Estimation of Soybean Yield Based on Physiological and Ecological Parameter and Meteorological Factor in China, *J Indian Soc Remote Sens*, 49, 873–886, <https://doi.org/10.1007/s12524-020-01269-3>, 2021.
- Li, X., Chen, M., He, S., Xu, X., He, L., Wang, L., Gao, Y., Tang, F., Gong, T., Wang, W., Xu, M., Liu, C., Yu, L., Liu, W., and Yang, W.: Estimation of soybean yield based on high-throughput phenotyping and machine learning, *Front. Plant Sci.*, 15, 1395760, <https://doi.org/10.3389/fpls.2024.1395760>, 2024.
- Liu, X. and Herbert, S. J.: Fifteen years of research examining cultivation of continuous soybean in northeast China: A review, *Field Crops Research*, 79, 1–7, [https://doi.org/10.1016/S0378-4290\(02\)00042-4](https://doi.org/10.1016/S0378-4290(02)00042-4), 2002.
- 680 Liu, X., Jin, J., Herbert, S. J., Zhang, Q., and Wang, G.: Yield components, dry matter, LAI and LAD of soybeans in Northeast China, *Field Crops Research*, 93, 85–93, <https://doi.org/10.1016/j.fcr.2004.09.005>, 2005.

- Liu, X., Jin, J., Wang, G., and Herbert, S. J.: Soybean yield physiology and development of high-yielding practices in Northeast China, *Field Crops Research*, 105, 157–171, <https://doi.org/10.1016/j.fcr.2007.09.003>, 2008.
- 685 Mei, Q., Zhang, Z., Han, J., Song, J., Dong, J., Wu, H., Xu, J., and Tao, F.: ChinaSoyArea10m: a dataset of soybean-planting areas with a spatial resolution of 10 m across China from 2017 to 2021, *Earth Syst. Sci. Data*, 16, 3213–3231, <https://doi.org/10.5194/essd-16-3213-2024>, 2024.
- Misaal, M. A., Zahra, S. M., Rasul, F., Imran, M., Noor, R., and Fahad, M.: Influence of Climate Change on Crop Yield and Sustainable Agriculture, in: *Climate Change Impacts on Natural Resources, Ecosystems and Agricultural Systems*, edited by: Pande, C. B., Moharir, K. N., Singh, S. K., Pham, Q. B., and Elbeltagi, A., Springer International Publishing, Cham, 209–223, https://doi.org/10.1007/978-3-031-19059-9_7, 2023.
- 690 Muhuri, A., Goita, K., Magagi, R., and Wang, H.: Soil Moisture Retrieval During Crop Growth Cycle Using Satellite SAR Time Series, *IEEE J. Sel. Top. Appl. Earth Observations Remote Sensing*, 16, 9302–9319, <https://doi.org/10.1109/JSTARS.2023.3280181>, 2023.
- 695 National Soil Survey Office: *Soil Species of China*, China Agriculture Press, Beijing, 924 pp., 1995.
- Ntakos, G., Prikaziuk, E., Ten Den, T., Reidsma, P., Vilfan, N., Van Der Wal, T., and Van Der Tol, C.: Coupled WOFOST and SCOPE model for remote sensing-based crop growth simulations, *Computers and Electronics in Agriculture*, 225, 109238, <https://doi.org/10.1016/j.compag.2024.109238>, 2024.
- Pang, A., Chang, M. W. L., and Chen, Y.: Evaluation of Random Forests (RF) for Regional and Local-Scale Wheat Yield Prediction in Southeast Australia, *Sensors*, 22, 717, <https://doi.org/10.3390/s22030717>, 2022.
- 700 Pasqualotto, N., Delegido, J., Van Wittenberghe, S., Rinaldi, M., and Moreno, J.: Multi-Crop Green LAI Estimation with a New Simple Sentinel-2 LAI Index (SeLI), *Sensors*, 19, 904, <https://doi.org/10.3390/s19040904>, 2019.
- Peng, G. and Yili, Z.: Research on Forest Phenology Prediction Based on LSTM and GRU Model, *Journal of Resources and Ecology*, 14, <https://doi.org/10.5814/j.issn.1674-764x.2023.01.003>, 2022.
- 705 Pinke, Z. and Lövei, G. L.: Increasing temperature cuts back crop yields in Hungary over the last 90 years, *Global Change Biology*, 23, 5426–5435, <https://doi.org/10.1111/gcb.13808>, 2017.
- Pu, L., Zhang, S., Yang, J., Chang, L., and Bai, S.: Spatio-Temporal Dynamics of Maize Potential Yield and Yield Gaps in Northeast China from 1990 to 2015, *IJERPH*, 16, 1211, <https://doi.org/10.3390/ijerph16071211>, 2019.
- Qiao, C., Cheng, C., and Ali, T.: How climate change and international trade will shape the future global soybean security pattern, *Journal of Cleaner Production*, 422, 138603, <https://doi.org/10.1016/j.jclepro.2023.138603>, 2023.
- 710 Qu, H., Li, X., Zhu, H., Wang, L., Qu, B., Wang, Q., Lv, J., Ji, Y., and Jiang, L.: Effects of combination of low temperature and excessive precipitation at seedling stage on soybean yield in high-latitude cold region, *Chinese Journal of Ecology*, 1–10, 2023.
- Ren, P., Li, H., Han, S., Chen, R., Yang, G., Yang, H., Feng, H., and Zhao, C.: Estimation of Soybean Yield by Combining Maturity Group Information and Unmanned Aerial Vehicle Multi-Sensor Data Using Machine Learning, *Remote Sensing*, 15, 4286, <https://doi.org/10.3390/rs15174286>, 2023a.
- 715

- Ren, Y., Li, Q., Du, X., Zhang, Y., Wang, H., Shi, G., and Wei, M.: Analysis of Corn Yield Prediction Potential at Various Growth Phases Using a Process-Based Model and Deep Learning, *Plants*, 12, 446, <https://doi.org/10.3390/plants12030446>, 2023b.
- 720 Shi, X. Z., Yu, D. S., Warner, E. D., Pan, X. Z., Petersen, G. W., Gong, Z. G., and Weindorf, D. C.: Soil Database of 1:1,000,000 Digital Soil Survey and Reference System of the Chinese Genetic Soil Classification System, *Soil Horizons*, 45, 129, <https://doi.org/10.2136/sh2004.4.0129>, 2004.
- Song, X.-P., Li, H., Potapov, P., and Hansen, M. C.: Annual 30 m soybean yield mapping in Brazil using long-term satellite observations, climate data and machine learning, *Agricultural and Forest Meteorology*, 326, 109186, 725 <https://doi.org/10.1016/j.agrformet.2022.109186>, 2022.
- Steduto, P., Hsiao, T. C., Raes, D., and Fereres, E.: AquaCrop—The FAO Crop Model to Simulate Yield Response to Water: I. Concepts and Underlying Principles, *Agron.j.*, 101, 426–437, <https://doi.org/10.2134/agronj2008.0139s>, 2009.
- Sun, X., Li, Q., Qiao, Y., Hu, Z., Zhang, X., and Liu, Y.: Warming and Drought in Hailun of Heilongjiang: Effects on Growth and Development of Soybean, *Chinese Agricultural Science Bulletin*, 38, 27–33, 2022.
- 730 Tan, J., Yang, P., Liu, Z., Wu, W., Zhang, L., Li, Z., You, L., Tang, H., and Li, Z.: Spatio-temporal dynamics of maize cropping system in Northeast China between 1980 and 2010 by using spatial production allocation model, *J. Geogr. Sci.*, 24, 397–410, <https://doi.org/10.1007/s11442-014-1096-0>, 2014.
- Tian, H., Wang, P., Tansey, K., Zhang, J., Zhang, S., and Li, H.: An LSTM neural network for improving wheat yield estimates by integrating remote sensing data and meteorological data in the Guanzhong Plain, PR China, *Agricultural and Forest Meteorology*, 310, 108629, <https://doi.org/10.1016/j.agrformet.2021.108629>, 2021. 735
- Urda, C., Rezi, R., Varga, A. G., Negrea, A., Muntean, E., Sopterean, L., and Duda, M. M.: EXPLORING THE IMPACT OF SOWING DATES ON SOYBEAN YIELD, SEED QUALITY AND TRYPSIN INHIBITOR ACTIVITY, *AGROLIFE SCIENTIFIC JOURNAL*, 13, 223–230, 2024.
- Von Bloh, M., Nóia Júnior, R. D. S., Wangerpohl, X., Saltık, A. O., Haller, V., Kaiser, L., and Asseng, S.: Machine learning for soybean yield forecasting in Brazil, *Agricultural and Forest Meteorology*, 341, 109670, 740 <https://doi.org/10.1016/j.agrformet.2023.109670>, 2023.
- Wang, B., Chen, C., Liu, D., Asseng, S., Yu, Q., and Yang, X.: Effects of climate trends and variability on wheat yield variability in eastern Australia, *Clim. Res.*, 64, 173–186, <https://doi.org/10.3354/cr01307>, 2015.
- Wang, C., Linderholm, H. W., Song, Y., Wang, F., Liu, Y., Tian, J., Xu, J., Song, Y., and Ren, G.: Impacts of Drought on Maize and Soybean Production in Northeast China During the Past Five Decades, *IJERPH*, 17, 2459, 745 <https://doi.org/10.3390/ijerph17072459>, 2020.
- Wang, X., Zhu, L., Hao, Y., Wang, Z., Xue, L., Ding, K., and Huang, X.: Impacts of aerosol meteorological feedback on China's yield potential of soybean, *Meteorological Applications*, 31, e2198, <https://doi.org/10.1002/met.2198>, 2024.
- 750 Xie, Q., Dash, J., Huete, A., Jiang, A., Yin, G., Ding, Y., Peng, D., Hall, C. C., Brown, L., Shi, Y., Ye, H., Dong, Y., and Huang, W.: Retrieval of crop biophysical parameters from Sentinel-2 remote sensing imagery, *International Journal of Applied Earth Observation and Geoinformation*, 80, 187–195, <https://doi.org/10.1016/j.jag.2019.04.019>, 2019.
- Xie, Y. and Huang, J.: Integration of a Crop Growth Model and Deep Learning Methods to Improve Satellite-Based Yield Estimation of Winter Wheat in Henan Province, China, *Remote Sensing*, 13, 4372, <https://doi.org/10.3390/rs13214372>, 2021.

- 755 Xu, J., Du, X., Dong, T., Li, Q., Zhang, Y., Wang, H., Xiao, J., Zhang, J., Shen, Y., and Dong, Y.: NortheastChinaSoybeanYield20m: an annual soybean yield dataset at 20 m in Northeast China from 2019 to 2023, <https://doi.org/10.5281/ZENODO.14263102>, 2024.
- Yang, S., Hu, L., Wu, H., Ren, H., Qiao, H., Li, P., and Fan, W.: Integration of Crop Growth Model and Random Forest for Winter Wheat Yield Estimation From UAV Hyperspectral Imagery, *IEEE J. Sel. Top. Appl. Earth Observations Remote Sensing*, 14, 6253–6269, <https://doi.org/10.1109/JSTARS.2021.3089203>, 2021.
- 760 Yildirim, T., Moriiasi, D. N., Starks, P. J., and Chakraborty, D.: Using Artificial Neural Network (ANN) for Short-Range Prediction of Cotton Yield in Data-Scarce Regions, *Agronomy*, 12, 828, <https://doi.org/10.3390/agronomy12040828>, 2022.
- Yu, Q., You, L., Wood-Sichra, U., Ru, Y., Joglekar, A. K. B., Fritz, S., Xiong, W., Lu, M., Wu, W., and Yang, P.: A cultivated planet in 2010 – Part 2: The global gridded agricultural-production maps, *Earth Syst. Sci. Data*, 12, 3545–3572, <https://doi.org/10.5194/essd-12-3545-2020>, 2020.
- 765 Zhang, Y., Liu, M., Kong, L., Peng, T., Xie, D., Zhang, L., Tian, L., and Zou, X.: Temporal Characteristics of Stress Signals Using GRU Algorithm for Heavy Metal Detection in Rice Based on Sentinel-2 Images, *IJERPH*, 19, 2567, <https://doi.org/10.3390/ijerph19052567>, 2022.
- 770 Zhao, G., Wang, J., Fan, W., and Ying, T.: Vegetation net primary productivity in Northeast China in 2000-2008: Simulation and seasonal change, *Ying yong sheng tai xue bao = The journal of applied ecology / Zhongguo sheng tai xue xue hui, Zhongguo ke xue yuan Shenyang ying yong sheng tai yan jiu suo zhu ban*, 22, 621–30, 2011.
- Zhao, J., Wang, C., Shi, X., Bo, X., Li, S., Shang, M., Chen, F., and Chu, Q.: Modeling climatically suitable areas for soybean and their shifts across China, *Agricultural Systems*, 192, 103205, <https://doi.org/10.1016/j.agsy.2021.103205>, 2021.
- Zhao, J., Wang, Y., Zhao, M., Wang, K., Li, S., Gao, Z., Shi, X., and Chu, Q.: Prospects for soybean production increase by closing yield gaps in the Northeast Farming Region, China, *Field Crops Research*, 293, 108843, <https://doi.org/10.1016/j.fcr.2023.108843>, 2023a.
- 775 Zhao, L., Li, Q., Chang, Q., Shang, J., Du, X., Liu, J., and Dong, T.: In-season crop type identification using optimal feature knowledge graph, *ISPRS Journal of Photogrammetry and Remote Sensing*, 194, 250–266, <https://doi.org/10.1016/j.isprsjprs.2022.10.017>, 2022.
- 780 Zhao, Y., Han, S., Zheng, J., Xue, H., Li, Z., Meng, Y., Li, X., Yang, X., Li, Z., Cai, S., and Yang, G.: ChinaWheatYield30m: a 30 m annual winter wheat yield dataset from 2016 to 2021 in China, *Earth Syst. Sci. Data*, 15, 4047–4063, <https://doi.org/10.5194/essd-15-4047-2023>, 2023b.
- Zheng, L. and Zhang, X.: Harvest time monitoring data of Shengyang Station in Liaoning Province from 1998 to 2008, National Ecosystem Science Data Center, <https://doi.org/10.12199/nesdc.ecodb.mon.2020.dp2011.sya.004.>, 2021.
- 785 Zhuo, W., Fang, S., Gao, X., Wang, L., Wu, D., Fu, S., Wu, Q., and Huang, J.: Crop yield prediction using MODIS LAI, TIGGE weather forecasts and WOFOST model: A case study for winter wheat in Hebei, China during 2009–2013, *International Journal of Applied Earth Observation and Geoinformation*, 106, 102668, <https://doi.org/10.1016/j.jag.2021.102668>, 2022.
- 790 Zhuo, W., Huang, H., Gao, X., Li, X., and Huang, J.: An Improved Approach of Winter Wheat Yield Estimation by Jointly Assimilating Remotely Sensed Leaf Area Index and Soil Moisture into the WOFOST Model, *Remote Sensing*, 15, 1825, <https://doi.org/10.3390/rs15071825>, 2023.