

Invariant Features for Global Crop Type Classification

Xin-Yi Tong¹, Sherrie Wang^{1,2,3}

¹ *Laboratory for Information and Decision Systems, MIT*

² *Department of Mechanical Engineering, MIT*

³ *Institute for Data, Systems, and Society, MIT*

Email: {xytong, sherwang}@mit.edu

Abstract

Accurately obtaining crop type and its spatial distribution at a global scale is of critical importance for food security, agricultural policy-making, and sustainable development. The rise of remote sensing technology offers an efficient solution for large-scale crop classification. However, the limited availability of reliable ground samples in most regions constrains the applicability of remote sensing-based crop classification methods across geographic areas. To address the performance decline of crop classification under geospatial shifts, this study focuses on identifying remote sensing features that are invariant to geographic variation and proposes strategies to enhance their cross-regional generalization. We constructed a global crop type dataset, CropGlobe, containing 300,000 pixel-level samples from eight countries across five continents and covering six major food and industrial crops: corn, soybeans, rice, wheat, sugarcane, and cotton. By its broad geographical coverage, CropGlobe enables a systematic evaluation of generalization under three transfer scenarios: cross-country, cross-continent, and cross-hemisphere. We compare the transferability of temporal multi-spectral features (Sentinel-2-based 1D/2D median features and harmonic coefficients) and hyperspectral features (from the Earth Surface Mineral Dust Source Investigation). To improve the generalization of features under spectral and phenological shifts, we design CropNet, a lightweight and robust convolutional neural network tailored for pixel-level crop classification, coupled with data augmentation strategies, including time shift, time scale, and magnitude warping, which simulate realistic variations in crop phenology across regions. Experimental results show that the 2D median features derived from Sentinel-2 consistently exhibit the strongest invariance across all transfer scenarios. Data augmentation further improves robustness, particularly under conditions of insufficient diversity in training data. This work identifies feature representations with greater invariance that enhance the geographic transferability of crop type classification, and suggests a promising paradigm toward scalable, low-cost applications across globally diverse regions. Our data and code will be available publicly at <https://x-ytong.github.io/project/CropGlobe.html>.

1 Introduction

Accurately acquiring crop type and its spatial distribution information at a global scale is of critical importance for ensuring food security, guiding agricultural production, and formulating agricultural policies [1–3]. With the ongoing challenges posed by climate change, population growth, and rising food demand, dynamic monitoring of agricultural systems has become a core issue for global sustainable development. The spatial and temporal patterns of crop types not only influence regional food supply and economic development but also have profound impacts on global ecosystem services, carbon cycles, and water resource regulation [4, 5].

Traditionally, crop type information has been obtained through field surveys and agricultural censuses, which are often labor-intensive and costly, hindering the timely updating of agricultural data [6]. In recent years, benefiting from the development of Earth observation technologies and machine learning methods, remote sensing has played an increasingly prominent role in agricultural monitoring. For example, the

Cropland Data Layer (CDL) [7] from the U.S. Department of Agriculture, the Annual Crop Inventory (ACI) [8] from Agriculture and Agri-Food Canada, and the Crop Map of England (CROME) [9] from the UK Rural Payments Agency are all high-quality remote sensing crop type mapping products. Nevertheless, these products are still mainly concentrated in countries or regions with strong statistical capacities, and it remains difficult to obtain such reliable crop type information in most parts of the world. The fundamental reason is that when no ground truth data is available in the target region, machine learning models often suffer from drastic performance degradation when applied to previously unseen samples [10, 11]. As a result, many countries or regions can only access crop area estimates at the provincial or national level, which are often of low quality.

One potential solution is to invest more resources into collecting ground truth data in the target areas, or to seek alternative non-traditional data sources, such as using crowdsourced labels instead of survey-based labels [12, 13], or having agricultural experts collect reference labels from Google Street View (GSV) [14, 15]. However, these methods are still costly, time-consuming, and suffer from inconsistent and unverifiable annotation quality. Another promising and complementary solution is to extract remote sensing features that are invariant to geographic variation. If crop types exhibit sufficient commonality across different regions, then existing models can be directly transferred to new areas without relying on ground labels, thus enabling cross-region crop type classification.

Currently, remote sensing-based crop classification mainly relies on two key types of features. The first is temporal multi-spectral features, which represent phenological characteristics of crops throughout the entire growing season [16, 17]. Physiological changes during different growth stages are reflected in the spectral responses, forming temporal patterns that can serve as a strong basis for crop identification. For instance, reflectance data from Sentinel-2 (S2) or Landsat is structured as a time series, and temporal features such as harmonic coefficients [18, 19] or biweekly/monthly median values [20, 21] can be extracted to represent crop growth dynamics for large-scale crop type mapping. The second is hyperspectral features, which capture the biophysical attributes of crops at a specific time point. Hyperspectral sensors such as the Earth Sensing Imaging Spectrometer (DESI) [22] and the PRecurSore IperSpettrale della Missione Applicativa (PRISMA) [23] provide hundreds of contiguous narrow bands that can characterize pigments, moisture, nitrogen content, and structural differences among crops [24, 25]. These indicators have strong discriminative ability, especially for fine-grained crop classification involving spectrally similar species. In addition, synthetic aperture radar (SAR) and light detection and ranging (LiDAR) have also been used to identify certain crop types. SAR, such as Sentinel-1 (S1), is insensitive to weather, which can compensate for the cloud-induced data gaps in optical imagery [26, 27]. Its ability to capture crop structural differences can assist in identifying certain crop types, such as rice [28]. LiDAR, such as the Global Ecosystem Dynamics Investigation (GEDI), provides information on canopy height and vertical structure, helping to distinguish between short and tall crops (e.g., soybean vs. maize) [29, 30]. However, radar-based data suffer from inherent drawbacks: SAR data are prone to geometric distortions and speckle noise, whereas LiDAR observations remain sparse in space and time, making them unsuitable as primary features for general-purpose, global-scale crop classification.

Although phenological and biophysical features each have their own advantages in crop type classification, their transferability is constrained by geographic variation. For temporal multi-spectral features, differences in climate, planting dates, management practices, and growth stages across regions cause changes in the phenological patterns of the same crop. For example, when biweekly median S2 features from the Northern Hemisphere are used to classify soybeans in the Southern Hemisphere, the performance drops sharply compared to intra-hemispheric transfers [31]. This may be because temporal multi-spectral features only capture the overall seasonal trends, lacking the ability to learn fine-scale contextual changes at specific time points. To address this issue, Temporal Convolutional Neural Networks (TempCNN) [32] provide a potential solution. TempCNN restructures the temporal multi-spectral features into a 2D format (spectral \times time matrix) and uses a convolutional neural network (CNN) to jointly capture spectral and temporal context, which may help enhance the geographic invariance of crop features. However, existing research on 2D temporal features remains limited to local-scale studies, lacking systematic validation at the global level. Recent advances

in transformer-based architectures, such as Presto [33], have demonstrated strong capability in performing highly generalizable spatiotemporal modeling for remote sensing data through self-supervised learning. This transferability typically comes at the cost of large model size, heavy computational demand, and reliance on extensive unlabeled pretraining data in addition to labeled datasets.

In terms of hyperspectral features, differences in soil background, atmospheric conditions, and planting density may cause shifts in crop spectral responses, resulting in inconsistent spectral signatures for the same crop across regions. Although the Earth Observing-1 (EO-1) Hyperion hyperspectral data have been used to construct agricultural crop spectral libraries in Central Asia [34] and the United States [35], where it achieves promising classification performance in each respective region, no prior work has explored the geographic invariance of hyperspectral features at a global scale.

These observations point to a notable research gap in understanding the geographic invariance of remote sensing features. In response, this study aims to address two key questions: (1) Which remote sensing features exhibit invariance to geographic variation for crop type classification at the global scale? (2) Under conditions where phenological and spectral profiles vary due to environmental factors; how can the invariance of features be further enhanced? To this end, we constructed a globally distributed crop type dataset, referred to as CropGlobe, by combining publicly available crop type products with S2 multi-spectral time-series and Earth Surface Mineral Dust Source Investigation (EMIT) hyperspectral data. CropGlobe contains 300,000 pixel-level samples from eight countries (Argentina, Australia, Belgium, China, France, the Netherlands, the United Kingdom, and the United States). We focus on six major grain and economic crops: corn, soybeans, rice, wheat, sugarcane, and cotton, which are widely cultivated and economically significant, forming the foundation of global agricultural systems.

On the CropGlobe dataset, we define three cross-region transfer scenarios: cross-country, cross-continent, and cross-hemisphere, and systematically compare temporal multi-spectral features (including 1D/2D median and harmonic coefficients) with hyperspectral features. To cope with the low-dimensional nature of pixel-level features, we propose a lightweight yet highly generalizable CNN model, CropNet. By incorporating spatial dropout and limiting downsampling depth, CropNet effectively prevents overfitting and enhances feature robustness. To further address the phenological misalignment and spectral variation across regions, we introduce a data augmentation strategy for temporal multi-spectral features, including time shift, time scale and magnitude warping transformations, with the goal of improving feature invariance against phenological and spectral shifts.

Our contributions are summarized as follows:

- We constructed a global crop type dataset, CropGlobe, spanning five continents and eight countries, focusing on six key global crops. The dataset is characterized by high quality and rich sample diversity, providing a benchmark resource for algorithm development and comprehensive analysis in crop type classification.
- We design a CNN model, CropNet, customized for pixel-level crop type classification, which outperforms widely used CNN baselines, including ResNet50 [36], EfficientNetV2-S [37], and ConvNeXt-Tiny [38], while maintaining low parameter complexity. In addition, we introduce a tailored data augmentation strategy for temporal multi-spectral features, which improves classification stability and yields up to a 7% accuracy gain under conditions of limited training diversity, offering a practical solution rarely addressed in previous studies.
- We compare different remote sensing features under diverse geographic transfer settings. To the best of our knowledge, this is the first systematic evaluation of the transferability of temporal multi-spectral and hyperspectral features for cross-country, cross-continent, and even cross-hemisphere crop type classification. We find that the 2D median features of S2 exhibit the strongest invariance globally. In the most challenging cross-hemisphere scenario, the overall classification accuracy achieves above 85%. This result demonstrates the practical feasibility of accurate global crop type mapping with substantially reduced local reference data requirements.

2 CropGlobe Dataset

To facilitate our investigation into feature invariance in global crop type classification, we constructed the CropGlobe dataset using 2023 as the reference year. It is built upon publicly available crop type reference products collected from eight countries across five continents: Argentina, Australia, Belgium, China, France, the Netherlands, the United Kingdom, and the United States. Two types of satellite data are incorporated: S2 multi-spectral time-series and EMIT hyperspectral imagery. We selected S2 for its relatively high spatial resolution and frequent revisit cycle, which are advantageous for crop monitoring. Meanwhile, the launch of the EMIT mission has enabled hyperspectral imagery with unprecedented spatial and temporal coverage, allowing the study of cross-geographical invariance in hyperspectral features.

We focus on six globally significant food and industrial crops: corn, soybeans, rice, wheat, sugarcane, and cotton, whose geographical distribution is displayed in Fig. 1. These crops are central to global agriculture, trade, industry, and food security. Including an “Other” class, CropGlobe consists of seven crop categories and contains more than 300,000 samples, each corresponding to a georeferenced pixel-level location on the Earth’s surface.

The following sections provide detailed descriptions of the satellite data sources (Section 2.1), crop type products (Section 2.2), and sample extraction procedures (Section 2.3).

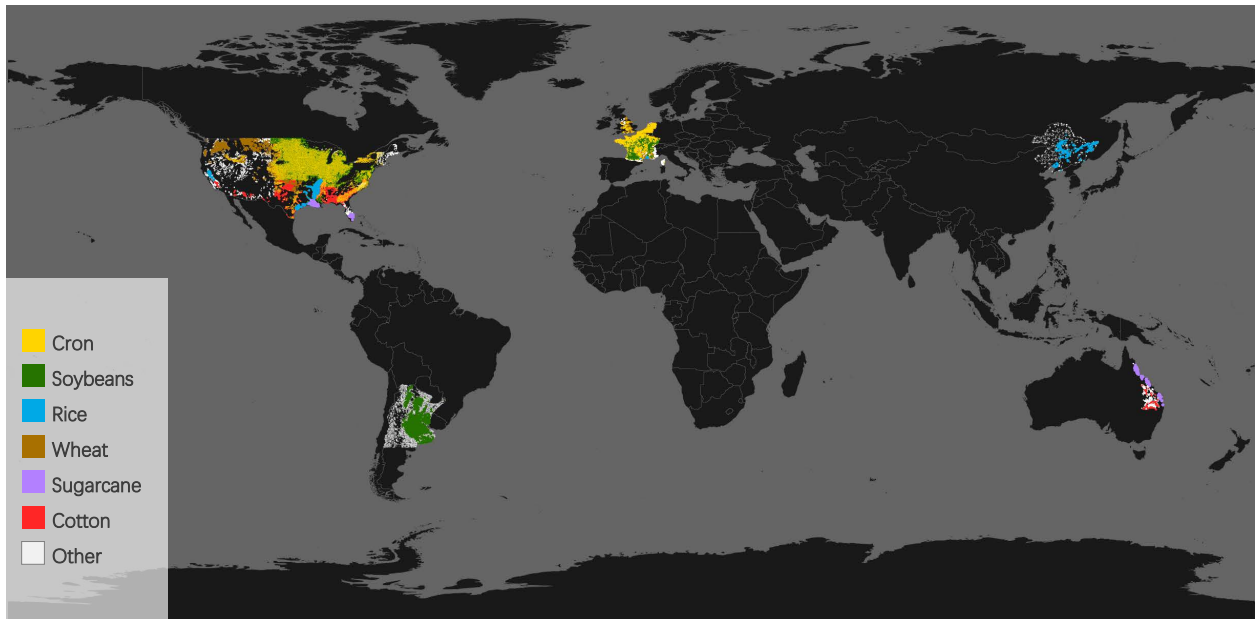


Figure 1: The geographical distribution and category system of the CropGlobe dataset. It contain 300,000 pixel-level samples from eight countries across five continents: Argentina, Australia, Belgium, China, France, the Netherlands, the United Kingdom, and the United States.

2.1 Satellite Data

Sentinel-2: S2 is part of the Copernicus Programme led by the European Space Agency (ESA), with a mission focused on the dynamic monitoring of terrestrial surfaces and vegetation. The system consists of two satellites: Sentinel-2A, launched in 2015, and Sentinel-2B, launched in 2017. Both satellites carry the Multi-Spectral Instrument (MSI), which acquires imagery in 13 spectral bands. Bands with a 10-meter spatial resolution include B2: 490 nm, Blue; B3: 560 nm, Green; B4: 665 nm, Red; and B8: 842 nm, Near Infrared

(NIR). Bands with a 20-meter resolution include B5: 705 nm, Red Edge 1; B6: 740 nm, Red Edge 2; B7: 783 nm, Red Edge 3; B8A: 865 nm, Narrow NIR; B11: 1610 nm, Shortwave Infrared 1 (SWIR 1); and B12: 2190 nm, Shortwave Infrared 2 (SWIR 2). Bands with a 60-meter resolution are B1: 443 nm, Coastal Aerosol; B9: 945 nm, Water Vapor; and B10: 1375 nm, Cirrus. Combined, the two satellites provide a global revisit frequency of five days. With the relatively high spatial resolution, S2 is well suited for multi-spectral time-series analysis to capture phenological patterns and differentiate crop types. It has become one of the most widely used data sources for agricultural dataset development and global crop mapping, such as CropHarvest [39] and WorldCereal [40].

Earth Surface Mineral Dust Source Investigation: EMIT is a hyperspectral imaging mission developed by the Jet Propulsion Laboratory of the National Aeronautics and Space Administration (NASA) and launched on July 14, 2022 (although preliminary data became accessible in August 2022, comprehensive spatial coverage was only established in 2023). The mission’s core objective is to monitor mineral dust on the surface of arid and semi-arid regions in order to understand its influence on the Earth system.

EMIT employs imaging spectroscopy to capture narrow, contiguous spectral bands spanning wavelengths from the visible to the shortwave infrared (380–2500 nm), with a spectral resolution of approximately 7.4 nm and a spatial resolution of around 60 meters. Each image pixel is associated with 285 spectral bands, of which 242 are retained after excluding those affected by atmospheric absorption windows. Onboard the International Space Station (ISS), EMIT has an approximate revisit cycle of three days, with the exact interval varying according to the orbital parameters of the station and Earth’s rotation.

Although EMIT is optimized for arid zones, it still provides partial coverage over agricultural regions worldwide. Its full spectral coverage encompasses all wavelengths relevant to crop type identification, enabling its use in evaluating the invariance of hyperspectral features at the global scale.

2.2 Reference Data

The following crop type mapping products serve as our reference data sources, chosen based on their availability, quality, and the geographic diversity of their coverage across continents.

Argentina (ARG) — South America Soybean (SAS): SAS [41] is a binary classification map developed to monitor soybean expansion across South America between 2000 and 2019, and has since been extended to 2023. It has a spatial resolution of 30 meters and focuses on assessing the indirect impact of soybean expansion on deforestation. The maps are generated through decision tree applied to multi-spectral time-series imagery, with model training based on field-verified samples. The reported overall accuracy is approximately 94–96%. As SAS only contains a binary distinction (soybean vs. non-soybean), the “Other” class in our dataset may in practice include corn, rice, wheat, sugarcane, or cotton. Data access: <https://glad.umd.edu/projects/commodity-crop-mapping-and-monitoring-south-america>.

Australia (AUS) — Queensland Seasonal Crop (QSC): QSC [42] is a seasonal crop classification product used to monitor large-scale cropping systems in Queensland, with the original data provided in vector format. Since 1988, multi-spectral time-series imagery has been used to generate two seasonal maps per year—one for winter crops (June–October) and another for summer crops (November–May). The 2023 summer crop map used in this study includes cotton, sugarcane, and an “Other” class that may contain corn, soybeans, and additional crop types. Reported accuracies are high, with user accuracy around 95% and producer accuracy about 90% for major crops [43]. Data access: <https://www.qld.gov.au/environment/land/management/mapping/statewide-monitoring/crops>.

Belgium (BEL) — Landbouwegebruikspercelen (LGP): LGP [44] is an annual crop declaration dataset maintained by the Belgian Department of Agriculture and Fisheries since 2008. Crop types are assigned to individual field parcels based on European standards, through either government surveys or farmer self-reporting. The original dataset contains 297 detailed categories in polygon vector format, which we processed to retain relevant classes and converted into raster format. Data access: <https://landbouwcijfers.vlaanderen.be/open-geodata-landbouwegebruikspercelen>.

China (CHN) — Crop Type in Northeast China (CTNC): CTNC [45] provides annual 10-meter

crop type maps for Northeast China from 2017 to 2019, with coverage later extended to 2023 by the original authors. The dataset is derived from S2 time-series imagery using a crop-specific ontology-based information fusion knowledge graph (OIFKG), and achieves producer accuracy greater than 92%. As only rice is included in this study, other crops such as corn, soybeans, wheat, sugarcane, and cotton are grouped under the “Other” class. Data access: https://figshare.com/articles/dataset/Mapping_crop_type_in_Northeast_China_during_2017-2023_Including_code_for_classification_/25346038.

France (FRA) — Registre Parcellaire Graphique (RPG): RPG [46] is a national agricultural parcel database maintained by the French Agency for Services and Payment (ASP), based on farmer-submitted declarations of field boundaries and crop types, and has been published annually since 2007. Although fields without ASP subsidies are not included, the RPG captures about 98% of French agricultural land. The original data, provided in polygon vector format with 372 crop categories, were simplified to include selected classes and converted to raster format. Data access: <https://geoservices.ign.fr/rpg>.

United Kingdom (GBR) — Crop Map of England (CROME): CROME [9] is an annual crop classification product maintained by the UK Department for Environment, Food and Rural Affairs (Defra), with updates available from 2016 to 2023. It is automatically generated using random forest classification applied to combined S1 radar and Planet optical imagery, and is verified using field survey and visual inspection. The original dataset, comprising 81 categories, is published in polygon vector format with hexagonal units. Reported accuracies for major crop classes are relatively high, with both user and producer accuracies exceeding 83%, whereas some non-major classes may have producer accuracies below 60%. We converted it into raster format and retained only the crop classes of interest. Data access: <https://environment.data.gov.uk/dataset/a27312b5-d6c9-4710-ad5e-382d727c1b05>.

Netherlands (NLD) — Basisregistratie Gewaspercelen (BRP): BRP [47] has been published annually by the Netherlands Enterprise Agency (Rijksdienst voor Ondernemend Nederland, RVO) since 2009. Agricultural parcel boundaries are determined based on the Dutch Agricultural Area Register, with farmers delineating their fields and reporting the crop types. The original dataset, provided in polygon vector format, contains 371 detailed categories. We converted it into raster format and selected the crop classes of interest. Data access: https://service.pdok.nl/rvo/brpgewaspercelen/atom/v1_0/basisregistratie_gewaspercelen-brp.xml.

United States (USA) — Cropland Data Layer (CDL): CDL [7] is a crop type mapping product released by the U.S. Department of Agriculture (USDA) National Agricultural Statistics Service (NASS), updated annually since 2008 (with some states since 1997). The maps are generated using decision tree classifiers trained on field survey data and satellite observations from Landsat 8/9 and S2 collected throughout the growing season. Each 30-meter pixel is assigned a crop type code that corresponds to a specific crop category. The original data contain 254 crop categories. Reported accuracies for major crop categories are typically in the range of 85–95% for both producer and user accuracies. For this study, we reorganized the data and selected the crop classes of interest, collecting data only from states where the reported accuracy meets the required standards (see Section 2.3 for details). Data access: <https://croplandcros.scinet.usda.gov/>.

2.3 Sample Collection

We collected samples for CropGlobe on the Google Earth Engine (GEE) platform through three main stages: satellite data preprocessing, reference data preprocessing, and grid-based spatial sampling. The procedures are described below.

Satellite Image Preprocessing: For multi-spectral time-series data, we used S2 Level-2A imagery (accessible via the GEE collection ‘COPERNICUS/S2_SR_HARMONIZED’), which provides surface reflectance corrected for atmospheric effects and cloud contamination. Cloudy pixels were filtered using the Cloud Score+ product (accessible via the GEE collection ‘GOOGLE/CLOUD_SCORE_PLUS/V1/S2_HARMONIZED’), with a threshold of 0.75, meaning only pixels with a non-cloud probability $\geq 75\%$ were retained. The Cloud Score+ product, generated from harmonized S2 Level-1C imagery, allows effective identification of cloud-free observations and removal of cloud and shadow pixels from Level-2A data. For Northern Hemisphere countries, we collected all

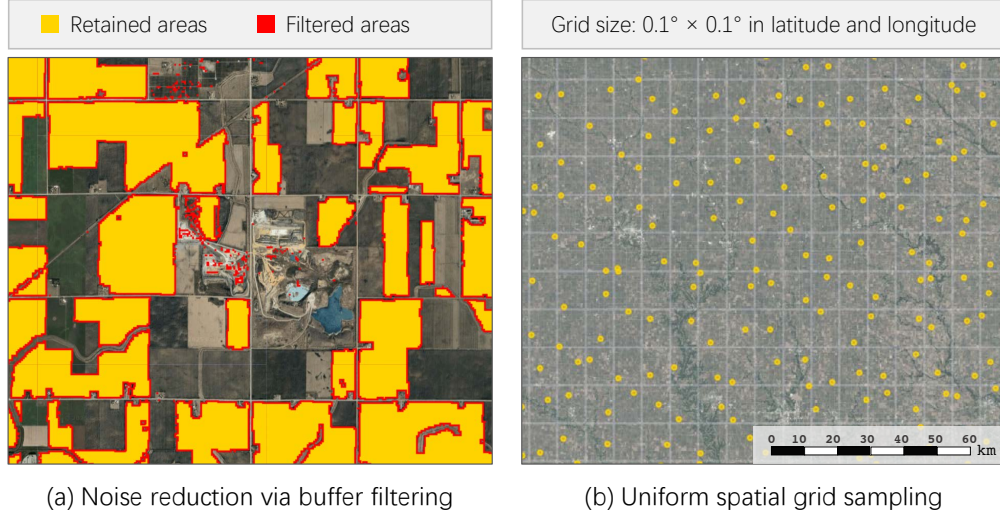


Figure 2: Illustration of noise reduction via buffered filtering and sample decorrelation through grid-based sampling.

available imagery from January 1 to December 31, 2023. For Southern Hemisphere countries, the period was shifted six months earlier, from July 1, 2022 to June 30, 2023, to account for the phenological offset between hemispheres.

For hyperspectral data, we used the EMIT Level-2A Estimated Surface Reflectance product (available via the GEE collection ‘NASA/EMIT/L2A/RFL’), which contains estimated surface reflectance, associated uncertainty, and cloud mask information. The built-in buffered cloud mask was used to exclude all pixels with cloud mask value equal to 1. To evaluate the utility of hyperspectral observations at different time points, we selected representative months that align with key crop growth stages, based on EMIT data availability over agricultural regions. For Northern Hemisphere countries, we used imagery from June and August 2023. We examined EMIT acquisitions throughout the main crop growing season and found that the overlap of cloud-free observations was largest between June and August. Including additional months such as May, July, or September resulted in negligible overlap and thus provided virtually no usable data. Median reflectance values were calculated for each month, and a combined availability mask—defined as the intersection of the two monthly cloud-free masks—was applied to retain valid pixels. For Southern Hemisphere countries, the timeframe was shifted by six months. As EMIT only established comprehensive spatial coverage in 2023, no November 2022 data were available over our study regions, and thus only February 2023 was used. In regions with data availability, monthly median hyperspectral reflectance was calculated using GEE.

Reference Data Preprocessing: Before importing the reference crop maps into GEE, we performed category extraction and merging. For BEL, FRA, and NLD, which are based on high-accuracy farmer declarations, we first removed non-agricultural classes such as buildings, water, wetlands, and forests, and then merged all categories outside of the six target crops (maize, soybean, rice, wheat, sugarcane, and cotton) into a single “Other” class.

For GBR and USA, where the maps are derived from automatic classification with available class-wise accuracy metadata, we applied strict filtering. In the GBR dataset, any class with user accuracy below 80% was removed (all six target classes exceeded this threshold), and all remaining classes were merged into the “Other” category. In the USA dataset, classification accuracy was evaluated at the state level; for each target class, only states with both producer and user accuracy above 85% were retained. In those states, non-agricultural classes were removed, and all classes indicating double cropping were excluded due to missing crop-specific timing information. Aside from the six target crops, only classes with accuracy above 80% were

Table 1: Class distribution of the CropGlobe dataset by country. Values marked with * correspond to samples derived from datasets with binary or ternary labels. The “Other” class may contain corn, soybeans, rice, wheat, sugarcane, or cotton, see Section 4.1 for processing details.

Country	Corn	Soybeans	Rice	Wheat	Sugarcane	Cotton	Other	Total
ARG	0	4,245	0	0	0	0	7,840*	12,085
AUS	0	0	0	0	5,098	4,825	8,938*	18,861
BEL	1,985	0	0	1,458	0	0	2,041	5,484
CHN	0	0	1,581	0	0	0	1,482*	3,063
FRA	24,906	21,097	2,294	27,688	0	0	32,714	108,699
GBR	4,624	0	0	3,927	0	0	4,712	13,263
NLD	2,388	0	0	2,135	0	0	2,795	7,318
USA	22,224	17,093	17,613	23,226	9,897	14,799	28,427	133,279
Total	56,127	42,435	21,488	58,434	14,995	19,624	88,949	302,052

Table 2: Class distribution of the CropGlobe_subset dataset. Values marked with * correspond to samples derived from datasets with binary or ternary labels. The “Other” class may contain corn, soybeans, rice, wheat, sugarcane, or cotton, see Section 4.1 for processing details.

Country	Corn	Soybeans	Rice	Wheat	Other	Total
ARG	0	3,143	0	0	3,205*	6,348
BEL	1,765	0	0	1,077	2,015	4,857
CHN	0	0	986	0	940*	1,926
FRA	5,511	324	2,910	7,863	8,221	24,829
GBR	133	0	0	1,383	1,313	2,829
NLD	1,554	0	0	1,011	2,243	4,808
USA	14,132	8,322	2,483	10,036	12,003	46,976
Total	23,095	11,789	6,379	21,370	29,940	92,573

retained and merged into the “Other” class. Due to differences in product accuracy and geographic coverage, the filtering criteria for GBR and USA are not identical. For GBR, we relied only on user accuracy because the overall accuracy of CROME is lower than that of CDL, and its producer accuracy is generally lower than its user accuracy. Applying stricter requirements would have eliminated nearly all non-target classes, leaving insufficient samples for the “Other” class. In contrast, for USA, CDL is of higher accuracy and covers a much larger national area. To avoid sample redundancy, we therefore applied stricter criteria.

For ARG, AUS, and CHN, the original datasets are binary or three-class products, and thus no further class selection or merging was required.

For datasets originally provided in vector format, namely QSC (AUS), LGP (BEL), RPG (FRA), CROME (GBR), and BRP (NLD), we converted them into raster format at a spatial resolution of 10 meters and reprojected them to the WGS84 coordinate reference system. To further improve label reliability, we applied a 60-meter inward buffer (based on EMIT resolution) to all labeled parcels after importing the crop maps into GEE, which helped eliminate small noisy patches. As shown in Fig. 2 (a), erroneous labels distributed within sample areas are removed after applying buffer filtering.

Grid-Based Sampling: We matched the buffered crop maps with the filtered satellite data (processed separately for S2 and EMIT) in both space and time. The intersection of valid crop labels and available cloud-free pixels defined the final sampling region, ensuring each pixel had both a class label and satellite data. To reduce sample correlation, the study region was divided into $0.1^\circ \times 0.1^\circ$ cells in geographic coordinates (WGS84), and random samples were drawn within each cell at a spatial resolution of 60 meters (i.e., ensuring

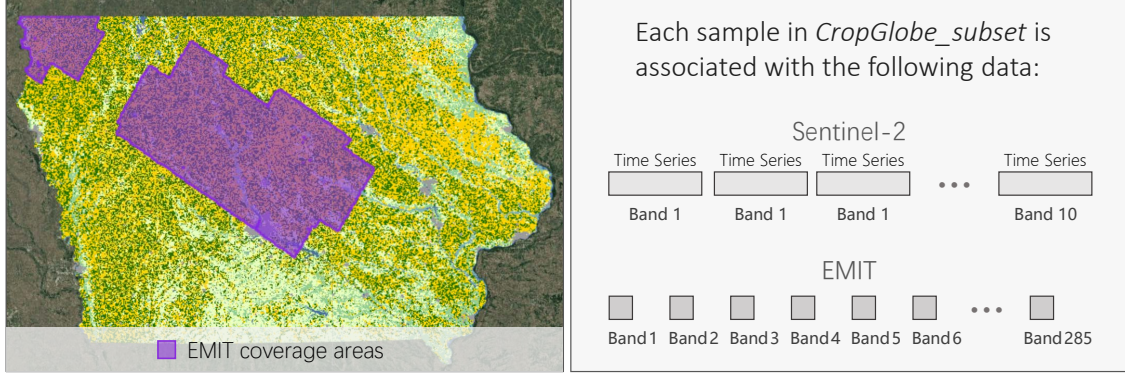


Figure 3: The coverage of EMIT is sparse. For example, in the U.S. state of Iowa, the overlapping area covered in both June and August 2023 is limited to the regions shown in purple. For each pixel within the purple region, both S2 time series and monthly median reflectance from EMIT are available.

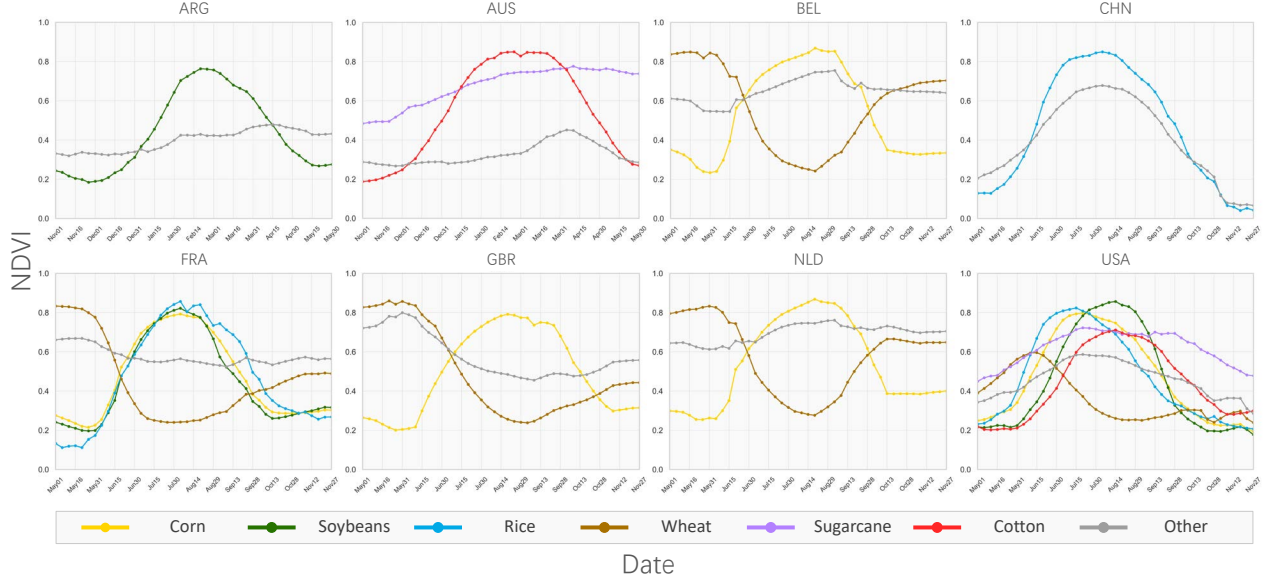


Figure 4: Average NDVI curves per country in the CropGlobe dataset. The temporal window spans May–November 2023 for Northern Hemisphere countries and November 2022–May 2023 for Southern Hemisphere countries, with 5-day intervals. The distinct temporal patterns across regions demonstrate the diversity of phenological dynamics in the dataset.

a minimum distance of 60 meters between any two sampled points), as shown in Fig. 2 (b). In USA, one sample was drawn per cell, while in other countries with smaller study areas, five samples were drawn per cell. Taking FRA as an example, a $0.1^\circ \times 0.1^\circ$ cell is approximately $7.7 \text{ km} \times 11.1 \text{ km}$, and five uniformly random samples with a minimum of 60-meter spacing remain spatially distinct.

For S2, we excluded bands B1, B9, and B10, and retained ten bands (B2, B3, B4, B8, B5, B6, B7, B8A, B11, B12). Each sample in the constructed CropGlobe dataset includes geolocation, a full-year time series of acquisition dates, and corresponding raw reflectance values for the ten spectral bands. The class distribution of the dataset is shown in Table 1. Due to the limited spatial coverage of EMIT (an illustrative example is

provided in Fig. 3), we derived a CropGlobe_subset dataset by selecting samples from CropGlobe that fall within EMIT-covered areas (i.e., June and August overlap in the Northern Hemisphere and February in the Southern Hemisphere). Each retained sample includes the same S2 multi-spectral time-series and 242-band EMIT median reflectance from single- or multi-date acquisitions, depending on regional data availability. The class distribution of CropGlobe_subset is shown in Table 2.

To illustrate the diversity and complexity of CropGlobe, we visualized the average Normalized Difference Vegetation Index (NDVI) [48] curves per country in Fig. 4. NDVI was calculated as $(\text{NIR} - \text{Red}) / (\text{NIR} + \text{Red})$, using B8 and B4 from S2. The temporal span was May to November 2023 for Northern Hemisphere countries and November 2022 to May 2023 for Southern Hemisphere countries, with a 5-day temporal window (although the full-year data were acquired, we found based on our subsequent sensitivity analysis that this period was the most representative for distinguishing crop types; details are provided in Section 4.5). NDVI curves showed clear differences in shape and magnitude across countries. For instance, corn exhibited different peak times and curve widths in different countries; wheat showed misaligned temporal patterns due to differences in planting and harvesting schedules. Peak NDVI values in FRA, GBR, and USA reached around 0.9, indicating dense vegetation, while lower values were observed in AUS and ARG, possibly reflecting less vigorous canopy growth. Even within the same class, curve smoothness, amplitude, and baseline noise varied significantly, highlighting the challenges for feature generalization. The “Other” class displayed particularly diverse curve patterns across countries, reflecting the heterogeneity of its composition.

3 Methodology

Achieving accurate crop type classification at a global scale presents a range of challenges, as the spectral and phenological characteristics of crops are influenced not only by environmental factors such as climate and soil, but also by the diversity of crop types and region-specific agricultural management practices [6, 10]. As a result, the same crop may exhibit significantly different spectral signatures and asynchronous phenological phases across regions. These variations exist across multiple scales: at the regional level, they lead to distributional shifts in the data, while at the local level, different crops may exhibit similar spectral signatures, increasing the risk of misclassification.

From a data structure perspective, pixel-level observations allow for dense temporal sampling but offer only a low-dimensional feature space per sample. When the features lack sufficient discriminative information, deep neural networks struggle to extract robust signals. In such cases, rather than capturing meaningful and transferable patterns, they tend to memorize noise or spurious correlations, ultimately undermining their ability to generalize to unseen data [49, 50].

To address these challenges, we design a lightweight yet effective CNN model named CropNet, tailored for learning robust features from pixel-level data. Furthermore, to improve the invariance of learned features across geographical domains, where both spectral responses and phenological phases may vary, we introduce a dedicated data augmentation strategy for multi-spectral time-series, aimed at enhancing robustness to temporal and spectral perturbations.

Our method comprises two main components: construction of features (Section 3.1 and Section 3.2), and construction of CropNet (Section 3.3).

3.1 Construction of Temporal Multi-Spectral Data

To identify features that are most invariant to geographic variation, we compare temporal multi-spectral features and hyperspectral features. For the temporal multi-spectral setting, we examine two types of representation: 1D/2D median features and harmonic coefficients. As the hyperspectral EMIT data lack full temporal coverage, we use its monthly median composites directly. Below, we detail the construction of the median and harmonic features.

1D/2D median features: To construct consistent median features from raw S2 observations, we apply

a windowed median resampling strategy. For each pixel, we define a time span of interest with a start day-of-year (DOY) T_s and an end DOY T_e , which typically covers the crop growing season. Although data are collected throughout the year, not all periods contribute equally to crop type classification. The selected time span allows the model to focus on the most relevant phenological phases.

This time span is divided into equal-length temporal windows of size d , resulting in $t = (T_e - T_s)/d$ discrete bins. For each bin, if valid (non-cloudy) observations are available, we compute the median reflectance across all dates within the bin. If no valid observations are found due to cloud coverage, the missing value is filled using linear interpolation along the temporal axis. If the time series contains too many gaps to allow reliable interpolation, the sample is discarded.

Each pixel is ultimately represented as a one-dimensional feature vector of length $10 \times t$, where 10 corresponds to the selected S2 bands, and t is the number of time bins.

Specifically, for constructing 2D features, this feature vector is then reshaped into a two-dimensional matrix of shape $10 \times t$, where rows represent spectral bands and columns correspond to temporal positions.

Harmonic coefficients: Harmonic coefficients are calculated based on four spectral bands: Narrow NIR, SWIR 1, SWIR 2, and NIR. In addition, the Green Chlorophyll Vegetation Index ($\text{GCVI} = \text{NIR} / (\text{Green} - 1)$) [51] is included. Frequency-domain features are extracted using third-order harmonic regression independently applied to each spectral band, resulting in a total of 35 features [14]. These four bands and the GCVI were selected following prior work [6], which demonstrated that using this combination achieves nearly the same crop-type classification performance as employing all bands and a broader set of vegetation indices.

3.2 Augmentation of Temporal Multi-Spectral Data

To enhance model robustness and generalization when using temporal multi-spectral features, we propose a data augmentation module specifically tailored for median-based time-series representations, incorporating three transformations: time shift, time scale, and magnitude warping, as illustrated in Fig. 5. These methods are designed to simulate natural variations in crop phenology, observation timing, and reflectance amplitude, which frequently occur across regions and seasons.

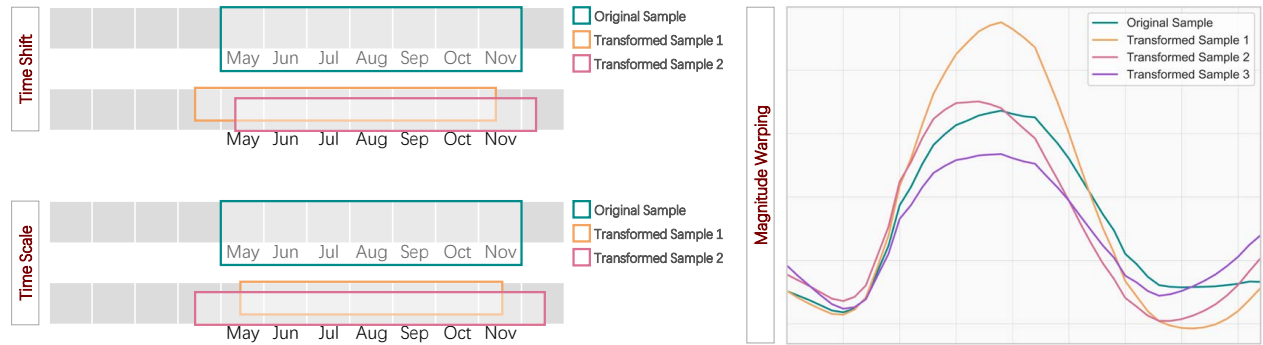


Figure 5: Illustration of three transformations for temporal data augmentation: time shift (randomly shifts a fixed temporal window forward or backward), time scale (compresses or stretches the time axis), and magnitude warping (smoothly distorts the amplitude over time).

Time Shift simulates phenological shifts by randomly moving the entire time window forward or backward by 0 to 10 days (see Section 4.5 for details). Given the original day-of-year (DOY) range $[T_s, T_e]$, a random offset is sampled and applied to both start and end DOYs, resulting in a shifted range $[T'_s, T'_e]$. Median values are then recalculated over this new interval using the same number of time bins as the original, ensuring feature length consistency. This augmentation accounts for planting date differences or climatic variations across regions.

Time Scale adjusts the temporal extent of the sampling window to mimic variability in crop development duration. The start and end DOYs are independently perturbed by values sampled from the range $[-30, +10]$ (see Section 4.5 for details), potentially expanding or contracting the time span by up to 40 days. To maintain a fixed input dimension, the temporal window is recomputed such that the total number of bins remains unchanged. This method captures differences in growing season length due to crop varieties, management practices or climate.

Magnitude Warping perturbs the shape of each spectral band’s temporal curve by applying a smooth, multiplicative modulation, aiming to mimic real-world deviations such as sensor noise, sub-pixel heterogeneity, or biomass variability (e.g., the same crop may even have different biomass in different fields) [52]. For median features of each sample, a random warping function is generated by cubic spline interpolation over a few anchor points (knots), whose values are drawn from a Gaussian distribution centered at 1.0. This function is then applied to scale the temporal profile of each spectral band. Based on empirical analysis (Section 4.5), we set the standard deviation of the Gaussian distribution to 0.2 and the number of anchor points (knots) to 5. To avoid over-regularization, each sample undergoes the same warping transformation across all spectral bands. This augmentation helps the model learn more invariant and generalizable representations under varying growing conditions.

All three augmentations operate on the raw multi-spectral time-series. Although the exact timing of crop growth is often a strong discriminative feature for crop type classification within a specific region, reliance on such region-specific cues limits cross-regional generalization. In contrast, our proposed transformations encourage the model to capture more broadly transferable features, thereby enhancing the model’s robustness to spectral and temporal variability and enabling more stable performance across regions.

3.3 Classification Model: CropNet

We design a lightweight CNN, named CropNet, specifically for pixel-level classification of crop type. It performs convolutions jointly along the spectral and temporal dimensions to capture dependencies across both domains, as presented in Fig. 6. CropNet consists of four convolutional blocks, each containing two 3×3 convolutional layers with batch normalization and ReLU (Rectified Linear Unit) activations. Down sampling is performed in the first and third blocks via stride-2 convolutions, retaining more temporal resolution while still allowing the model to capture multi-scale spectral-temporal patterns. After the final convolutional stage, global average pooling reduces the spatial dimensions to 1×1 , followed by a linear layer for classification. The output is a probability distribution over crop types.

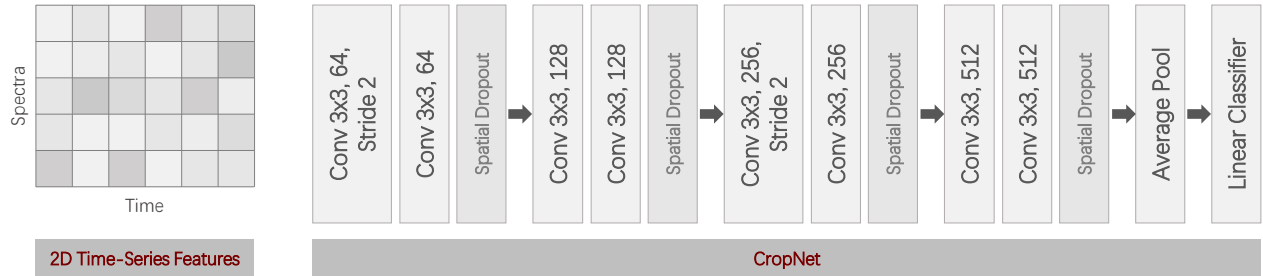


Figure 6: Architecture of CropNet. Downsampling is applied in the first and third blocks, enabling multi-scale feature extraction while preserving temporal resolution. Spatial dropout is employed to improve regularization by zeroing out entire feature channels.

To improve regularization and prevent overfitting on small per-pixel features, we apply spatial dropout [53] after each convolutional block. Unlike standard dropout, which randomly drops individual values, spatial dropout zeros out entire feature channels, forcing the network to learn more robust and distributed repre-

sentations across the spectrum-time space. This is particularly useful in our case, as the input tensor lacks contextual spatial structure and is more prone to overfitting due to limited feature diversity per sample.

We evaluated several architectural variants with different depths and channel widths. Our experiments showed that a four-stage, nine-layer configuration achieves the best trade-off between performance and model complexity. It provides stable accuracy across test regions and is a minimal design for optimal generalization. Deeper networks showed no significant gains and were more prone to overfitting, while shallower ones consistently underperformed.

For one-dimensional inputs or datasets structured as flat time-series vectors, such as harmonic coefficients or hyperspectral reflectance, the 2D convolutions in CropNet can be easily replaced with 1D convolutions. In this configuration, the network architecture remains structurally identical, treating the input as a single-channel signal and performing convolutions along the feature axis.

4 Experimental Results

We organize our experiments into three groups according to the geographic transfer scenarios: cross-country (FRA \rightarrow BEL, FRA \rightarrow NLD, FRA \rightarrow GBR), cross-continent (USA \rightarrow FRA, USA \rightarrow CHN, FRA \rightarrow USA), and cross-hemisphere (USA \rightarrow AUS, USA \rightarrow ARG, FRA \rightarrow ARG). Here, “ \rightarrow ” indicates that the model is trained on the source region and directly applied to the target region without accessing any data from the target domain during training. Under these settings, we conduct comprehensive comparisons and analyses to investigate feature invariance across different conditions.

On the CropGlobe dataset, we compare 1D and 2D median-based features against harmonic coefficient features derived from S2 time-series. On the CropGlobe_subset dataset, we compare temporal multi-spectral features with single-/multi-date hyperspectral features. In addition, we perform sensitivity analysis on the selection of temporal window and time span for constructing multi-spectral features, and conduct ablation studies to assess the effectiveness of the proposed data augmentation strategies. This section is structured as follows: (1) Experimental Setup (Section 4.1); (2) 1D/2D Median Features vs. Harmonic Coefficient (Section 4.2); (3) Temporal Multi-Spectral Features vs. Hyperspectral Features (Section 4.3); (4) Comparison of CNN Models (Section 4.4); (5) Sensitivity Analysis and Ablation Study (Section 4.5).

4.1 Experimental Setup

Class alignment: Since available crop categories differ across countries and the definition of the “Other” class also varies, we apply a label alignment strategy for fair evaluation. For AUS, ARG, and CHN [41, 42, 45], where the datasets are binary or ternary and the “Other” class may implicitly include target classes (e.g., soybeans and wheat may exist in AUS “Other”), we adopt the following rule: when using these datasets as test regions, any predicted class not explicitly defined within their available categories is reassigned to “Other” before computing accuracy. For example, samples predicted as soybean or wheat in AUS are reassigned to “Other” during evaluation. In the FRA \rightarrow USA experiment, since FRA contains no sugarcane or cotton samples, we remove these two classes from USA when conducting evaluation; in this case, they cannot simply be merged into “Other” because FRA entirely lacks these categories. For all other transfer scenarios, no special category alignment is applied, as the training regions cover all classes present in the test regions and the definition of “Other” remains consistent.

Feature processing: For temporal multi-spectral features, the temporal window used for computing medians is 5 days, and the time span for both median and harmonic coefficient ranges from May to November (see Section 4.5 for details). When applying data augmentation to median features, each of the three augmentation transformations—time shift, time scale, and magnitude warping—is applied independently with a probability of 0.5. For hyperspectral features, single-date data are directly used as raw vectors without additional processing. For multi-date hyperspectral features, two vectors from different dates are concatenated into a combined feature vector.

Model comparison: We adopt Random Forest (RF) with 500 trees as a baseline. In addition, we compare CropNet with several standard CNN architectures, including ResNet50 [36] (classical deep residual network with skip connections), EfficientNetV2-S [37] (compound-scaled lightweight model balancing depth, width, and resolution), and ConvNeXt-Tiny [38] (modernized convolutional architecture inspired by Transformer designs). To accommodate the low-dimensional input features, we adjust the degree of spatial downsampling to preserve sufficient resolution throughout these networks, keeping their overall architecture intact. All models are trained with a learning rate of 0.0001, batch size of 256, and for 50 epochs.

Evaluation metrics: We report Overall Accuracy (OA) and mean F1 score (mF1) to evaluate model performance. The mF1 score is computed as the macro-averaged F1 across all categories, reflecting the model’s ability to balance both overestimation and underestimation for each class. For each experiment, we fix different random seeds and repeat the training ten times, reporting the mean and standard deviation of results.

Table 3: Crop type classification performance on CropGlobe. Best results are highlighted. Each experiment is repeated ten times, and the mean and standard deviation are reported. “AUGM” denotes the use of data augmentation. Values in parentheses indicate performance gain from augmentation.

	Method	FRA → BEL		FRA → NLD		FRA → GBR	
		OA (%)	mF1 (%)	OA (%)	mF1 (%)	OA (%)	mF1 (%)
Cross-Country	RF_Harmonic	81.02 ± 0.16	83.15 ± 0.15	78.38 ± 0.22	80.02 ± 0.21	69.85 ± 0.14	71.35 ± 0.13
	RF_Median 1D	92.26 ± 0.11	94.55 ± 0.07	89.21 ± 0.07	90.47 ± 0.06	89.44 ± 0.14	90.16 ± 0.10
	CropNet_Harmonic	82.52 ± 0.53	86.23 ± 0.32	80.08 ± 0.87	83.87 ± 0.69	68.53 ± 1.70	71.53 ± 1.50
	CropNet_Median 1D	93.09 ± 0.69	94.39 ± 0.51	90.49 ± 0.38	91.27 ± 0.25	84.60 ± 0.89	85.19 ± 0.86
	CropNet_Median 2D	97.79 ± 0.30	97.96 ± 0.28	95.43 ± 0.44	95.55 ± 0.41	95.31 ± 0.91	95.35 ± 0.95
	CropNet_Median 2D_AUGM	98.34 ± 0.21 (+ 0.55)	98.47 ± 0.17 (+ 0.51)	95.82 ± 0.20 (+ 0.39)	95.88 ± 0.18 (+ 0.33)	95.80 ± 0.77 (+ 0.49)	95.88 ± 0.79 (+ 0.53)
	Method	USA → FRA		USA → CHN		FRA → USA	
		OA (%)	mF1 (%)	OA (%)	mF1 (%)	OA (%)	mF1 (%)
Cross-Continent	RF_Harmonic	53.42 ± 0.20	50.33 ± 0.22	71.79 ± 0.70	70.86 ± 0.80	45.09 ± 0.07	38.63 ± 0.10
	RF_Median 1D	75.34 ± 0.10	75.01 ± 0.12	72.95 ± 2.42	71.77 ± 2.82	57.22 ± 0.09	50.31 ± 0.10
	CropNet_Harmonic	53.28 ± 1.35	57.44 ± 1.20	66.32 ± 2.05	63.56 ± 2.71	53.24 ± 0.80	49.44 ± 1.01
	CropNet_Median 1D	66.74 ± 1.65	68.17 ± 1.61	60.52 ± 3.29	55.29 ± 4.96	51.14 ± 0.88	45.72 ± 1.32
	CropNet_Median 2D	86.90 ± 1.46	88.91 ± 1.25	93.17 ± 4.23	93.14 ± 4.31	70.40 ± 1.44	69.21 ± 1.57
	CropNet_Median 2D_AUGM	86.97 ± 0.90 (+ 0.07)	89.48 ± 0.60 (+ 0.57)	95.21 ± 0.38 (+ 2.04)	95.21 ± 0.38 (+ 2.07)	75.19 ± 1.54 (+ 4.79)	74.88 ± 1.98 (+ 5.67)
	Method	USA → AUS		USA → ARG		FRA → ARG	
		OA (%)	mF1 (%)	OA (%)	mF1 (%)	OA (%)	mF1 (%)
Cross-Hemisphere	RF_Harmonic	61.21 ± 0.34	58.14 ± 0.37	80.97 ± 0.11	75.98 ± 0.17	78.55 ± 0.12	73.32 ± 0.17
	RF_Median 1D	70.00 ± 0.48	66.36 ± 0.64	86.82 ± 0.16	84.29 ± 0.21	77.56 ± 0.41	70.03 ± 0.71
	CropNet_Harmonic	62.41 ± 1.67	63.11 ± 1.38	81.80 ± 0.96	77.18 ± 1.52	81.34 ± 1.01	77.83 ± 1.45
	CropNet_Median 1D	71.73 ± 1.65	71.58 ± 1.50	88.07 ± 0.69	85.86 ± 0.93	71.47 ± 2.02	58.32 ± 4.65
	CropNet_Median 2D	86.58 ± 1.55	86.61 ± 1.43	92.23 ± 0.82	91.31 ± 0.88	83.91 ± 3.77	79.81 ± 5.54
	CropNet_Median 2D_AUGM	87.16 ± 0.83 (+ 0.58)	87.28 ± 0.79 (+ 0.67)	92.45 ± 0.70 (+ 0.22)	91.60 ± 0.74 (+ 0.29)	91.04 ± 0.63 (+ 7.13)	89.75 ± 0.79 (+ 9.94)

4.2 2D Median Features Outperform 1D Median and Harmonic Coefficient

Table 3 presents the results of different feature–classifier combinations across cross-country, cross-continent, and cross-hemisphere scenarios. Overall, median features consistently outperform harmonic coefficients. Even when using a simpler classifier like RF, the performance of RF with 1D median features exceeds that of CropNet with harmonic features. In all settings, we observe that 2D median features significantly outperform their 1D counterparts, as 2D representations capture not only inter-band relationships but also contextual information along the temporal axis.

Data augmentation generally improves the stability of results, and brings notable performance gains in the FRA → USA and FRA → ARG experiments. We speculate that, although we adopted a higher

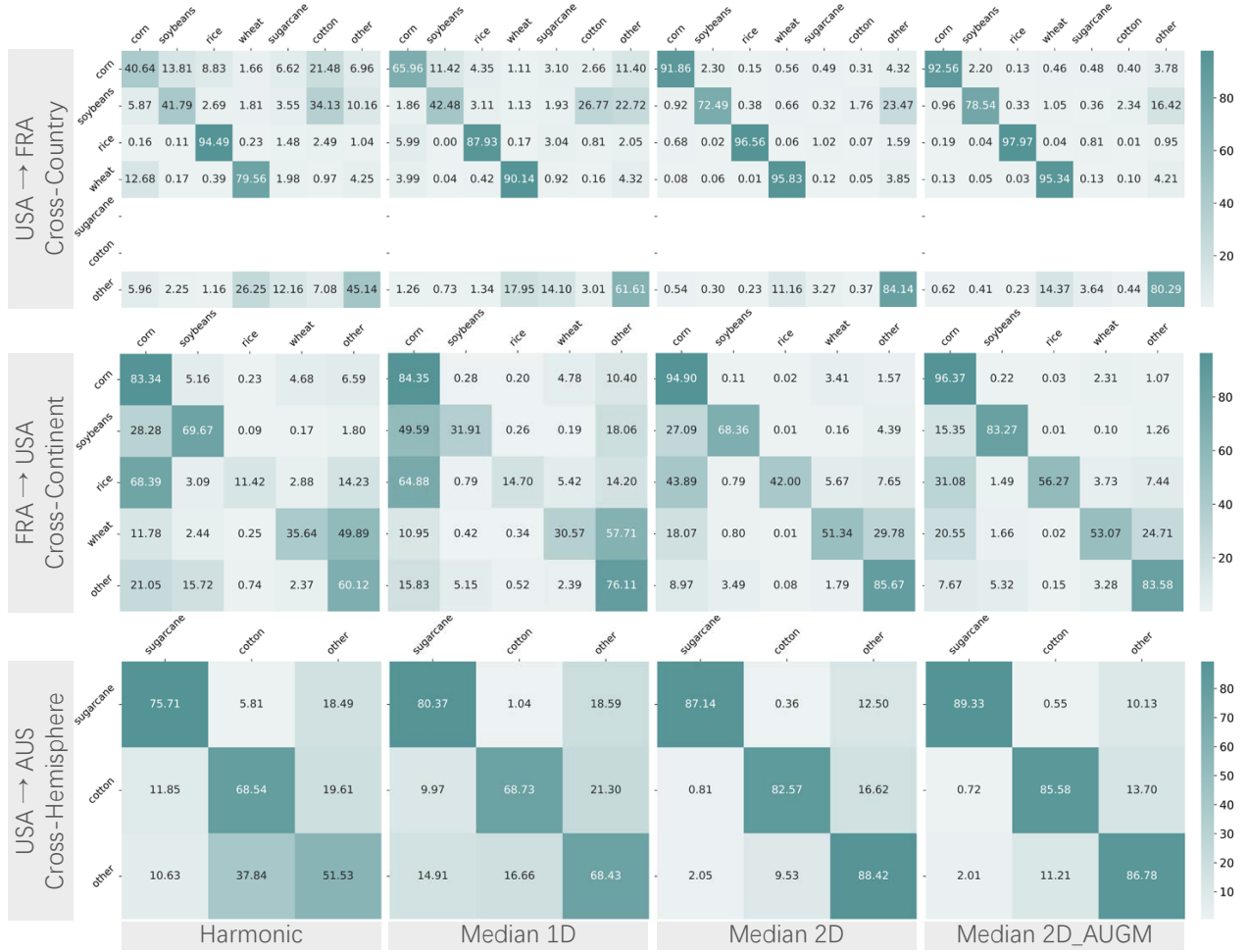


Figure 7: Confusion matrices for selected cross-region classification results from Table 3. 1D version of CropNet is used for harmonic and 1D median features, while 2D median features and augmented inputs are processed by 2D CropNet. Rows correspond to ground-truth crop types, and columns represent predicted labels. Higher values along the diagonal indicate better classification accuracy.

sampling density, the relatively small area of France still limits the diversity of training samples. When there is a large distribution shift between training and test samples, augmentation effectively compensates for insufficient data variety. In contrast, when training on USA data, the contribution of data augmentation is relatively limited. This is because USA covers a wide range of latitudes and longitudes, encompassing diverse crop-growing conditions across different climates and geographies, resulting in more representative training samples.

Classification accuracy is closely related to geographic distance. The best cross-country results generally exceed 95%, while cross-continent and cross-hemisphere results are slightly lower. In relatively simple binary classification settings (e.g., USA → CHN, USA → ARG, and FRA → ARG), cross-continent and cross-hemisphere accuracies can still exceed 90%. The most difficult case is FRA → USA, where the best result is 75.19%. This indicates that even with augmentation to simulate phenological variability, it is difficult to fully compensate for the lack of sample diversity in the training set.

Fig. 7 shows the confusion matrices for three cross-region experiments: USA → FRA, FRA → USA, and

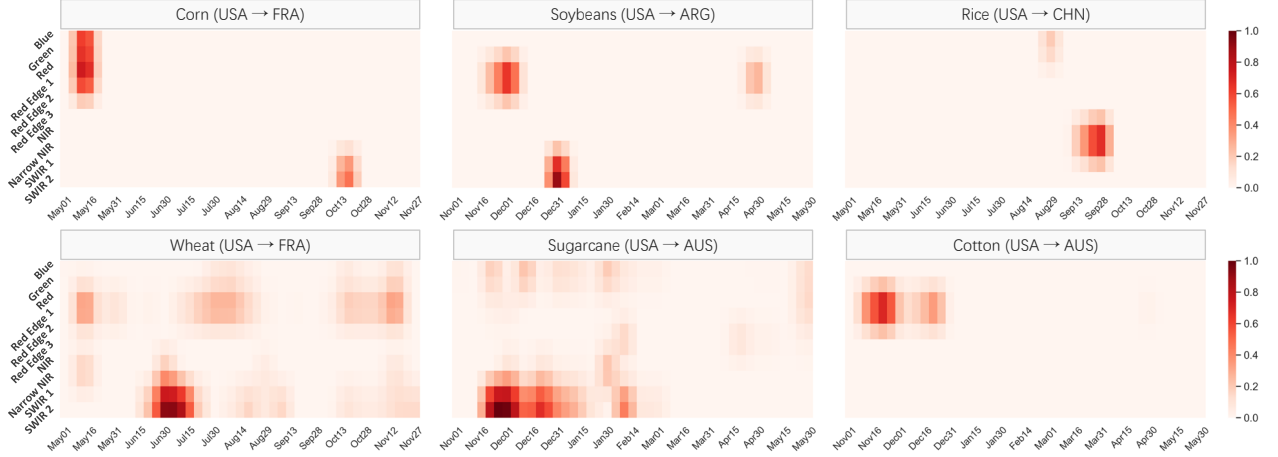


Figure 8: Visualization of 2D feature importance maps from 2D CropNet for six crop and region transfer settings. The x-axis indicates time, and the y-axis shows S2 spectral bands. For AUS and ARG, time has been shifted six months earlier. The highlighted regions indicate joint spectral-temporal areas that are most informative for classification.

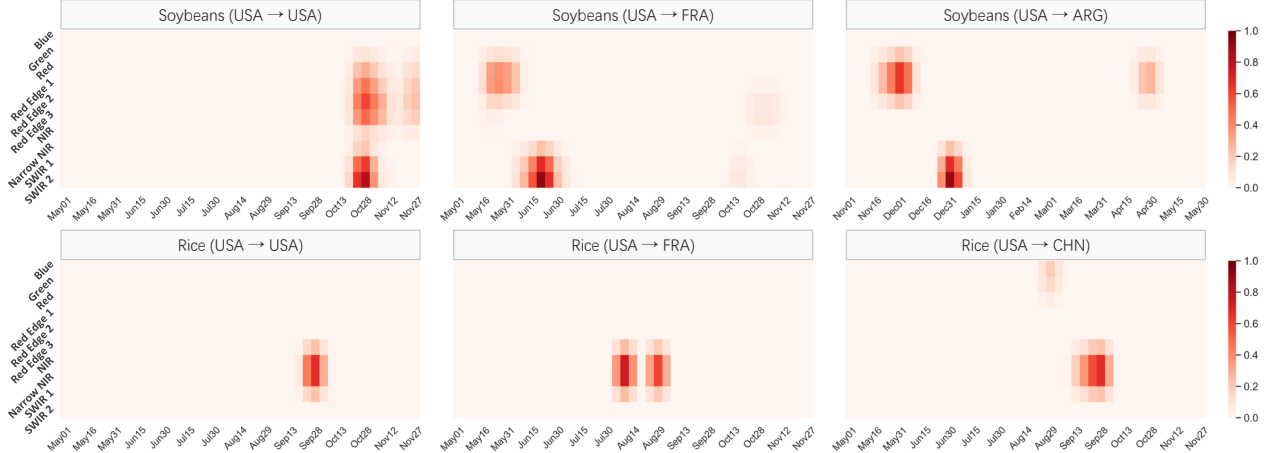


Figure 9: Feature importance maps of the same crop (soybeans and rice) in different target regions. The x-axis denotes time, and the y-axis indicates S2 spectral bands. While the time importance is distributed across different bins due to region-specific crop growth phenology, the most informative spectral bands remain consistent across regions.

USA → AUS, all using CropNet as classifier. Compared to 1D features, 2D median features consistently improve per-class accuracy, and data augmentation further enhances performance on top of that. Notably, in the cross-hemisphere setting of USA → AUS, all crop classes achieve over 85% accuracy. In contrast, the FRA → USA experiment yields less satisfactory results, particularly for wheat, due to the lower diversity of training data in FRA compared to USA. Nevertheless, data augmentation still leads to substantial improvements for soybeans and rice, highlighting its effectiveness even when sample diversity is limited.

We visualized the feature importance derived from 2D median using the CropNet_AUGM model. The importance was computed using Class Activation Maps (CAMs) [54], which estimate the contribution of each input region to the model’s prediction by backpropagating gradients from the predicted class through

intermediate convolutional layers. In this case, the CAMs highlight joint patterns across both the spectral and temporal dimensions.

Fig. 8 displays the importance maps for six crop types. Notably, the highlighted regions are not isolated points, but form continuous zones across both bands and time, indicating that the 2D features successfully capture spectral-temporal dependencies. This allows the model to learn more invariant and robust representations for crop classification under varying regional conditions. In general, the most informative spectral bands for crop identification are concentrated in Red, Red Edge 1, Red Edge 2, Narrow NIR, SWIR 1 and SWIR 2. These bands are known to be sensitive to chlorophyll content, leaf structure, and water content, which are important indicators of crop health and development.

It also can be seen that the temporal patterns vary across crops. For example, corn exhibits strong importance signals in early May over a relatively short period, corresponding to its rapid emergence and early growth phase. In contrast, wheat and sugarcane show broader temporal importance spans, reflecting their longer growing cycles and more extended phenological stages. This suggests that the model adapts to crop-specific developmental timelines when learning discriminative features.

To further investigate spatial generalization, we visualized CAMs for soybeans and rice under different target regions in Fig. 9. For a given crop, the important time periods vary significantly across regions, while the key spectral bands remain stable. For soybeans, the most informative bands are typically Red to Red Edge 2 and Narrow NIR to SWIR 2, which are sensitive to photosynthetic activity and moisture stress. For rice, the most critical bands are located in Red Edge 3 to SWIR 1 region, which are particularly responsive to canopy structure and flooding or water-related dynamics, consistent with the biophysical characteristics of paddy rice. These findings confirm that while crop phenology varies by environment, their underlying spectral responses remain consistent. By capturing spectral-temporal dependencies, the 2D model is able to extract shared patterns that generalize across regions [32].

To evaluate the intrinsic discriminative ability of different features and models without the added complexity of domain shifts, we conduct in-region classification experiments where training and testing samples are drawn from the same country. Table 4 summarizes the results where models are trained and tested within the same country. For each country, 80% of samples are randomly selected for training and 20% for testing, with experiments repeated ten times. The results show that median features consistently outperform harmonic coefficients, and 2D median features further outperform 1D median features. Data augmentation remains effective even under these in-region settings, improving both model stability and accuracy, despite limited domain shifts between training and testing data.

Table 4: In-region crop classification performance for FRA and USA. Each model is trained and tested within the same country using an 80%/20% train-test split, repeated ten times. ‘‘AUGM’’ denotes the use of data augmentation. It can be seen that data augmentation remains effective even in low domain shift settings.

Method	FRA → FRA		USA → USA	
	OA (%)	mF1 (%)	OA (%)	mF1 (%)
RF_Harmonic	84.55 ± 0.21	86.74 ± 0.21	79.32 ± 0.21	79.89 ± 0.22
RF_Median 1D	95.77 ± 0.15	96.41 ± 0.13	94.15 ± 0.14	94.64 ± 0.16
CropNet_Harmonic	88.74 ± 0.21	90.61 ± 0.17	88.14 ± 0.23	88.69 ± 0.24
CropNet_Median 1D	95.30 ± 0.13	96.14 ± 0.14	93.61 ± 0.18	94.14 ± 0.19
CropNet_Median 2D	98.44 ± 0.07	98.73 ± 0.06	98.16 ± 0.05	98.40 ± 0.04
CropNet_Median 2D_AUGM	99.37 ± 0.03	99.49 ± 0.02	99.28 ± 0.04	99.37 ± 0.04

It is noteworthy that the FRA dataset is based on farmer-reported data and thus highly reliable, while the USA dataset is automatically generated and its label accuracy is inherently imperfect. Despite this, our model achieves near 99% accuracy on the USA data, suggesting that some degree of overfitting to label noise may have occurred. However, the strong generalization observed in USA-to-other transfer experiments suggests that the learned features still generalize well despite noisy labels.

Another observation is that RF performs very well when sufficient labeled data are available within the

region of interest. This suggests that when enough high-quality samples can be collected locally, Random Forest classifiers deployed on cloud platforms such as GEE can serve as a practical solution for large-scale, high-accuracy crop type mapping. In future applications, it may be possible to first use 2D median features and CropNet to automatically generate reliable pseudo-labels for parts of the target region, which can then serve as training samples for Random Forest to scale up mapping in the same region.

Table 5: Crop type classification performance on CropGlobe_subset. Best results are highlighted. Each experiment is repeated ten times, and the mean and standard deviation are reported. All experiments are conducted using the CropNet model without data augmentation. The “S2_Median_2D” variant uses the 2D version of CropNet, while all other variants use the 1D version.

	Feature	FRA → BEL		FRA → NLD		FRA → GBR	
		OA (%)	mF1 (%)	OA (%)	mF1 (%)	OA (%)	mF1 (%)
Cross-Country	S2_Harmonic	79.45 ± 1.35	83.95 ± 0.80	77.99 ± 1.49	82.85 ± 0.97	68.23 ± 1.62	61.77 ± 1.85
	S2_Median 1D	93.85 ± 0.42	94.62 ± 0.31	93.00 ± 0.53	93.76 ± 0.40	89.96 ± 1.39	90.02 ± 1.01
	S2_Median 2D	97.56 ± 0.37	97.68 ± 0.33	97.19 ± 0.70	97.44 ± 0.62	94.84 ± 1.46	94.85 ± 1.15
	EMIT_Jun	80.53 ± 1.64	83.63 ± 0.73	79.90 ± 1.47	83.32 ± 0.77	90.76 ± 1.11	87.38 ± 1.57
	EMIT_Aug	82.25 ± 1.61	86.70 ± 1.14	78.73 ± 1.73	84.49 ± 0.93	73.58 ± 0.61	73.19 ± 0.72
	EMIT_Jun&Aug	91.41 ± 1.33	93.53 ± 1.14	90.15 ± 1.37	93.02 ± 0.86	90.60 ± 0.98	89.20 ± 1.07
	Feature	USA → FRA		USA → CHN		FRA → USA	
		OA (%)	mF1 (%)	OA (%)	mF1 (%)	OA (%)	mF1 (%)
Cross-Continent	S2_Harmonic	70.06 ± 0.78	61.47 ± 0.95	91.29 ± 0.89	91.27 ± 0.90	53.66 ± 1.69	55.22 ± 1.08
	S2_Median 1D	78.67 ± 0.96	67.91 ± 0.84	94.61 ± 0.44	94.61 ± 0.44	61.86 ± 1.72	57.56 ± 2.97
	S2_Median 2D	90.43 ± 0.94	83.63 ± 1.51	93.13 ± 0.97	93.12 ± 0.97	69.13 ± 1.86	68.20 ± 2.89
	EMIT_Jun	69.14 ± 1.17	55.23 ± 1.24	90.56 ± 0.49	90.52 ± 0.50	47.52 ± 0.74	47.74 ± 0.91
	EMIT_Aug	70.65 ± 0.75	60.18 ± 0.60	82.54 ± 0.38	82.19 ± 0.40	59.60 ± 0.87	53.89 ± 1.99
	EMIT_Jun&Aug	81.55 ± 0.91	68.66 ± 0.91	88.72 ± 1.03	88.65 ± 1.05	63.98 ± 1.36	62.51 ± 1.99
	Feature	USA → ARG		FRA → ARG			
		OA (%)	mF1 (%)	OA (%)	mF1 (%)		
Cross-Hemisphere	S2_Harmonic	78.07 ± 1.00	77.29 ± 1.17	68.23 ± 1.55	66.91 ± 2.00		
	S2_Median 1D	80.16 ± 3.59	79.27 ± 4.10	51.59 ± 0.59	36.12 ± 1.33		
	S2_Median 2D	94.89 ± 0.86	94.89 ± 0.86	53.00 ± 2.36	39.11 ± 4.82		
	EMIT_Feb	84.76 ± 2.02	84.36 ± 2.19	67.89 ± 2.58	64.57 ± 3.47		

4.3 Temporal Multi-Spectral Features Outperform Hyperspectral Features

Table 5 compares temporal multi-spectral features and hyperspectral features on the CropGlobe_subset dataset using the CropNet-1D/2D models. Some entries are missing due to data limitations, specifically, CropGlobe_subset lacks coverage in Australia and does not include Southern Hemisphere data for November 2022, rendering certain experiments infeasible.

The results show that hyperspectral features offer reasonable transferability, achieving over 80% accuracy in several settings (FRA → BEL, FRA → GBR, USA → CHN, and USA → ARG) using a single timestamp. When using data from both June and August, performance improves further, exceeding 90% in all cross-country scenarios. This is because hyperspectral data captures fine-grained spectral fingerprints of crops, including pigment composition, nitrogen levels, and biomass, which tend to exhibit some degree of invariance across geographic regions.

However, temporal multi-spectral features with sufficient time-series resolution consistently outperform hyperspectral features, particularly in challenging cases like USA → FRA and USA → ARG. These findings confirm the value of detailed temporal information in capturing crop phenology for classification.

Interestingly, in the FRA → ARG case, the harmonic coefficient achieves the best performance. This result contrasts with the findings on the full dataset reported in Table 3, where the best-performing features are Median 2D_AUGM. The discrepancy can be explained by differences in crop type distributions: in CropGlobe (see Table 1), there are sufficient soybean samples for FRA, whereas in CropGlobe_subset (see

Table 2), soybean samples are extremely limited, accounting for only 1.30% of all FRA samples. Since ARG involves a binary soybean classification, under such sample-limited conditions, low-dimensional features like harmonic coefficients may be less prone to overfitting and yield more robust results.

4.4 CropNet Outperforms Baseline CNN Models

Table 6 presents a comprehensive comparison between our proposed CropNet and three state-of-the-art CNN architectures, ResNet50, EfficientNetV2-S, and ConvNeXt-Tiny, on cross-region crop classification using 2D median features. No data augmentation is applied in this set of experiments. Notably, CropNet contains only 4.69 million parameters, substantially fewer than ResNet50 (23.57 M), EfficientNetV2-S (20.24 M), and ConvNeXt-Tiny (27.85 M), yet it consistently achieves competitive or superior performance across all transfer settings.

In the cross-country scenarios, CropNet achieves the highest OA and mF1 scores in all three cases. These results validate that CropNet effectively captures domain-invariant spectral-temporal structures while maintaining strong expressiveness, despite its lightweight nature.

In the cross-continent transfer experiments, which are inherently more challenging due to increased phenological and climatic discrepancies, CropNet maintains a clear advantage. It outperforms all baselines on USA \rightarrow FRA, FRA \rightarrow USA and matches or exceeds their performance on USA \rightarrow CHN.

In the cross-hemisphere setting, which introduces further challenges due to seasonal inversion and hemispheric phenology shifts, CropNet also achieves strong results. It outperforms larger models while maintaining compactness and efficiency. The results confirm that CropNet strikes an effective balance between model complexity and performance, offering strong generalization, compact architecture, and suitability for operational-scale applications such as continental or global crop-type mapping.

Table 6: Comparison of CNN models for crop type classification on CropGlobe. All models operate on 2D inputs and no data augmentation is applied. CropNet achieves competitive performance with significantly fewer parameters (4.69M), demonstrating its efficiency and effectiveness across all transfer settings.

CNN Model Parameters		ResNet50		EfficientNetV2-S		ConvNeXt-Tiny		CropNet (Ours)	
		23.57 M		20.24 M		27.85 M		4.69 M	
		OA (%)	mF1 (%)	OA (%)	mF1 (%)	OA (%)	mF1 (%)	OA (%)	mF1 (%)
Cross-Country	FRA \rightarrow BEL	96.89 \pm 0.30	97.16 \pm 0.24	97.35 \pm 0.44	97.57 \pm 0.41	97.34 \pm 0.92	97.61 \pm 0.77	97.79 \pm 0.30	97.96 \pm 0.28
	FRA \rightarrow NLD	93.27 \pm 0.57	93.39 \pm 0.60	94.38 \pm 0.41	94.47 \pm 0.39	94.00 \pm 0.83	94.23 \pm 0.71	95.43 \pm 0.44	95.55 \pm 0.41
	FRA \rightarrow GBR	92.56 \pm 0.97	92.67 \pm 1.01	93.20 \pm 1.23	93.34 \pm 1.20	93.75 \pm 1.18	93.82 \pm 1.22	95.31 \pm 0.91	95.35 \pm 0.95
Cross-Continent	USA \rightarrow FRA	81.43 \pm 1.74	81.78 \pm 2.17	83.70 \pm 1.32	85.03 \pm 1.59	80.93 \pm 1.37	83.04 \pm 1.53	86.90 \pm 1.46	88.91 \pm 1.25
	USA \rightarrow CHN	71.36 \pm 8.76	69.18 \pm 10.18	88.36 \pm 7.85	88.12 \pm 8.26	95.20 \pm 0.57	95.19 \pm 0.57	93.17 \pm 4.23	93.14 \pm 4.31
	FRA \rightarrow USA	59.87 \pm 1.34	55.26 \pm 1.78	63.95 \pm 1.93	59.99 \pm 2.57	61.15 \pm 1.66	56.71 \pm 1.60	70.40 \pm 1.44	69.21 \pm 1.57
Cross-Hemisphere	USA \rightarrow AUS	81.21 \pm 2.58	80.97 \pm 2.74	81.62 \pm 2.86	81.67 \pm 2.79	79.60 \pm 2.41	79.69 \pm 2.15	86.58 \pm 1.55	86.61 \pm 1.43
	USA \rightarrow ARG	88.25 \pm 0.93	86.31 \pm 1.18	89.02 \pm 1.15	87.47 \pm 1.42	89.39 \pm 0.81	87.92 \pm 0.97	92.23 \pm 0.82	91.31 \pm 0.88
	FRA \rightarrow ARG	80.67 \pm 2.67	75.31 \pm 4.23	80.90 \pm 2.46	75.74 \pm 3.71	80.94 \pm 2.27	75.60 \pm 3.69	83.91 \pm 3.77	79.81 \pm 5.54

4.5 Sensitivity Analysis and Ablation Study

We conduct a sensitivity analysis to investigate how different temporal window sizes and time spans affect model performance across transfer settings. As shown in Fig. 10, results are reported as mean OA averaged over all transfer directions within each setting (Cross-Country, Cross-Continent, and Cross-Hemisphere). The bar plots on the left evaluate different temporal windows, while the right plots assess different time spans.

In terms of temporal window, we observe that shorter windows (e.g., 5-day intervals) consistently lead to better classification results across all transfer scenarios. This suggests that longer intervals (e.g., 15 or 30 days) may compress multiple growth stages into a single observation, causing critical signals to be smoothed out. In contrast, shorter intervals preserve the full dynamics of crop development, capturing subtle growth rates and transition points rather than only the overall seasonal shape. As a result, the model can accurately learn phenological differences among crops. Regarding time span, the impact differs by transfer type. For

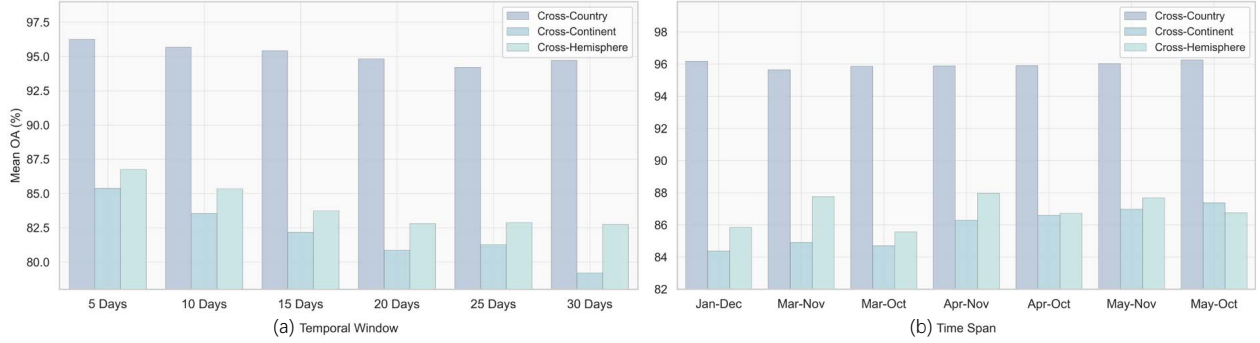


Figure 10: Impact of (a) different temporal window lengths and (b) different seasonal time spans on mean OA across Cross-Country, Cross-Continent, and Cross-Hemisphere settings.

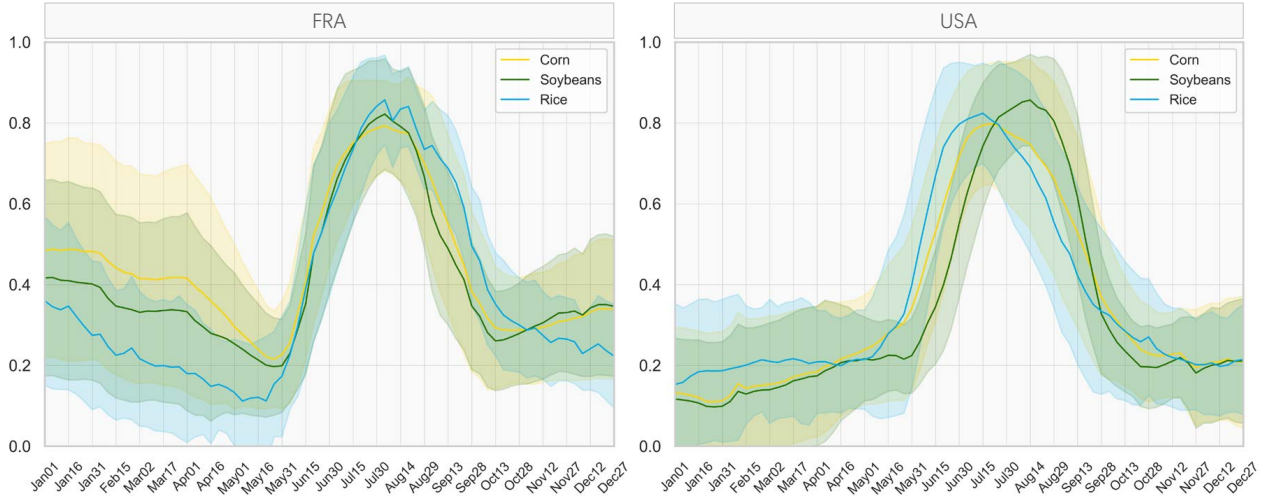


Figure 11: NDVI time series for corn, soybeans, and rice in FRA and USA over a full year. In France, NDVI values are already high from January to April, indicating the presence of vegetation during the early months of the year. This may be due to early-sown crops, winter crops, or crop rotation practices.

cross-country transfer, changing the time span has little effect. This is expected, as countries like FRA, BEL, NLD, and GBR are geographically close and share similar climate regimes, agricultural calendars, and crop phenology. Therefore, truncating or extending the time span does not substantially alter the informativeness of the temporal features.

In contrast, cross-continent results are more sensitive to time span changes. The best results are achieved using data from May to October, while including early-year months (January to March) leads to performance degradation. To analyze this, we visualize the average NDVI profiles of corn, soybeans, and rice in FRA and USA (Fig. 11). In FRA, we observe significant NDVI activity during January to March, particularly for corn and soybeans. This is likely due to crop rotation practices (e.g., planting rapeseed in the off-season), which elevate NDVI values and introduce noisy signals unrelated to the target summer crops. When such signals from non-growing season are included, the model may overfit to these phenological patterns present in the training region, which do not appear in target regions (e.g., USA, where early-year NDVI remains low). This mismatch leads to domain shifts and ultimately degrades generalization performance.

For cross-hemisphere transfer, the best performance is achieved using time spans from April to November.

Moreover, we find that including November consistently improves accuracy regardless of the start month. This is because crops such as wheat, sugarcane, and cotton often have longer growth cycles. In domain-shifted settings, time periods near the end of the season provide discriminative phenological cues that enhance cross-domain alignment. Considering the trade-offs across different settings, we select May to November as the final time span used in our main experiments, as it achieves consistently high performance in both cross-continent and cross-hemisphere scenarios.

We further conduct ablation study on data augmentation strategies, as presented in Table 7. The experiments are conducted under two challenging cross-region settings: FRA \rightarrow USA and FRA \rightarrow ARG. The results show that time shift and time scale individually bring only modest improvements over the baseline. However, the inclusion of magnitude warping on top of these two leads to a clear and substantial gain in both OA and F1 score across FRA \rightarrow USA and FRA \rightarrow ARG settings.

Table 7: Ablation study on the components of data augmentation. We evaluate the individual and combined impact of time shift, time scale, and magnitude warping on crop type classification in two cross-region settings: FRA \rightarrow USA and FRA \rightarrow ARG.

Components	FRA \rightarrow USA		FRA \rightarrow ARG	
	OA (%)	mF1 (%)	OA (%)	mF1 (%)
No Augmentation	70.40 \pm 1.44	69.21 \pm 1.57	83.91 \pm 3.77	79.81 \pm 5.54
Time Shift	71.23 \pm 1.52	70.45 \pm 1.69	85.21 \pm 1.55	81.90 \pm 1.86
Time Shift + Time Scale	72.31 \pm 1.58	71.72 \pm 1.78	87.42 \pm 1.02	84.17 \pm 1.43
Time Shift + Time Scale + Magnitude Warping	75.19 \pm 1.54	74.88 \pm 1.98	91.04 \pm 0.63	89.75 \pm 0.79

The superior performance of magnitude warping can be attributed to its ability to simulate realistic and continuous variations in spectral reflectance over time. Unlike time shift or time scale, which merely shift or stretch the entire temporal profile uniformly, magnitude warping introduces smooth, localized fluctuations that resemble real-world factors such as crop heterogeneity, sensor noise, or biophysical variability caused by weather, soil, and management differences. These subtle distortions help the model learn more flexible and robust temporal representations, allowing it to generalize better to unseen regions where crop development patterns may not follow a rigid shift or scale transformation. This also implies that the multi-spectral time series retains sufficient discriminative signals after perturbations, perhaps in the interrelationships among different spectral bands. These underlying relative relationships and their temporal co-evolution remain stable enough to be exploited by the model, ensuring that enough information is preserved for reliable classification.

To further refine our understanding of data augmentation strategies, we conduct a detailed sensitivity analysis of the parameters involved in time shift, time scale, and magnitude warping. These experiments are carried out under two challenging transfer settings: FRA \rightarrow USA and FRA \rightarrow ARG, and results are reported as the mean OA across both tasks.

Fig. 12 presents the results for individual application of time shift and time scale. In each case, only the corresponding augmentation method is used. For time shift, we observe that moderate perturbation—specifically the range of $[-10, +10]$ yields the best performance. This can be attributed to the fact that our base temporal span is already limited to May through November, as justified in earlier analyses (see Fig. 10). Including phenologically irrelevant months (e.g., December or early spring) introduces noise rather than useful variability, leading to reduced classification accuracy.

For time scale, where the start and end dates of the sampling window are perturbed independently, the best performance is achieved with a range of $[-30, +10]$. This suggests that stronger scaling transformations (particularly contracting the window) may help the model better adapt to variations in crop development duration due to climate or management differences. However, when the scaling becomes too extreme, performance declines slightly, likely due to the increased chance of including off-season or irrelevant observations that dilute the discriminative signal.

Fig. 13 investigates the effect of magnitude warping parameters. Here, only magnitude warping is applied

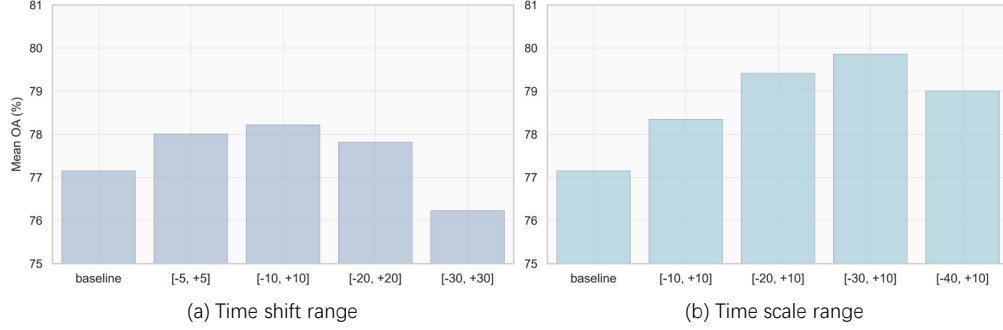


Figure 12: Sensitivity analysis of parameters for data augmentation: (a) Impact of time shift range on model performance; (b) Impact of time scale range. Only one transformation (either time shift or time scale) is applied in each experiment. Mean OA values are averaged over the FRA \rightarrow USA and FRA \rightarrow ARG transfer settings.

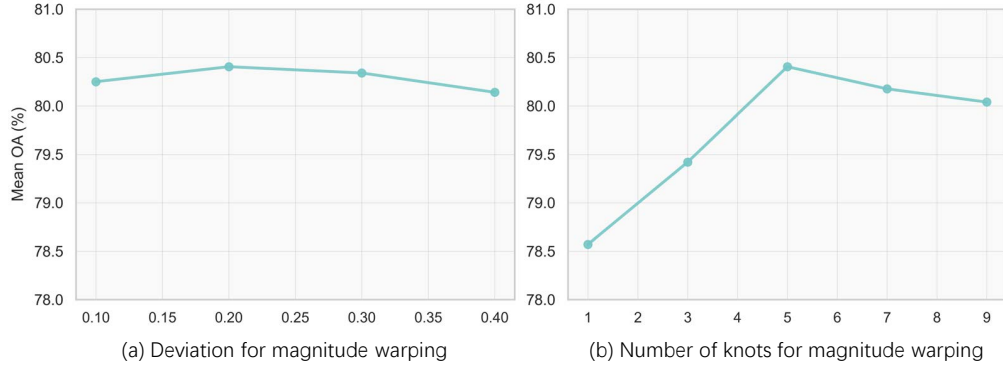


Figure 13: Sensitivity analysis of parameters for data augmentation: (a) Effect of Gaussian deviation for magnitude warping; (b) Effect of the number of spline knots for magnitude warping. Only magnitude warping is applied in this set of experiments. Mean OA values are averaged over the FRA \rightarrow USA and FRA \rightarrow ARG transfer settings.

in each experiment. The left plot shows the effect of varying the standard deviation of the Gaussian noise used to perturb the anchor points of the warping curve. While overall performance is relatively stable across deviation values, the best result is observed at 0.2, which introduces moderate amplitude variability without destabilizing the signal. Higher deviations do not consistently improve performance.

The right plot shows how the number of spline knots, which is used to define the smooth modulation curve, affects accuracy. Too few knots (e.g., 1 or 3) produce overly rigid distortions that do not effectively simulate natural fluctuations. In contrast, too many knots (e.g., 9) lead to excessive local variations, possibly introducing artificial artifacts that deviate from real-world crop behavior. The optimal number of knots is 5, which strikes a balance between smoothness and variability. This sensitivity analysis shows that moderate time shift, scale change, and magnitude distortion lead to better balance between realistic variation and stable performance.

5 Discussion

5.1 From Local to Global: Making Cross-Regional Crop Type Mapping Possible

Cross-regional crop type mapping has long been regarded as a formidable challenge due to natural differences in climate, phenology, and soil conditions, as well as the uneven availability of labeled data across regions. In this study, using the CropGlobe dataset collected from eight countries across five continents, we conducted a systematic evaluation of the transferability of different remote sensing features.

Our results reveal that S2 temporal multi-spectral features transfer better than commonly assumed. Even under substantial geographic shifts, a simple combination of S2 median features with a lightweight 2D convolutional model consistently delivers promising accuracy. Using the same median features, 2D CropNet show a clear gain over 1D CropNet. The likely reason is CNNs’ near-invariance to translations in the joint spectral–temporal plane. Along the temporal axis, when phenology shifts earlier or later (e.g., due to planting date or climate), the pattern largely translates and remains detectable. Along the spectral axis, small band shifts or amplitude changes (from sensor, atmosphere, cultivar, or management differences) do not erase local structures. Pooling and hierarchical feature aggregation further absorb such small deformations. By contrast, a 1D model that only scans along time struggles to capture cross-band local dependencies, limiting transferability.

Beyond this core finding, simple temporal data augmentation strategies, including time shift, time scale, and magnitude warping, can mitigate limited diversity in the training data to some extent. These strategies reduce the performance gap between data-rich and data-scarce regions, underscoring the potential of straightforward preprocessing to enhance model generalization.

To the best of our knowledge, this represents the first systematic investigation of transferability in crop type classification on a global scale. Our results show that combining off-the-shelf satellite imagery, strong feature representations, and lightweight architectures makes crop type mapping feasible across many regions worldwide with substantially fewer local labels. Sharing data or pretrained models across countries enable rapid mapping in new regions and lower the annotation cost for updates in regions already covered. We also expect these benefits to extend across years (e.g., under anomalous seasons), although we do not evaluate cross-temporal transfer here; we leave it for future work.

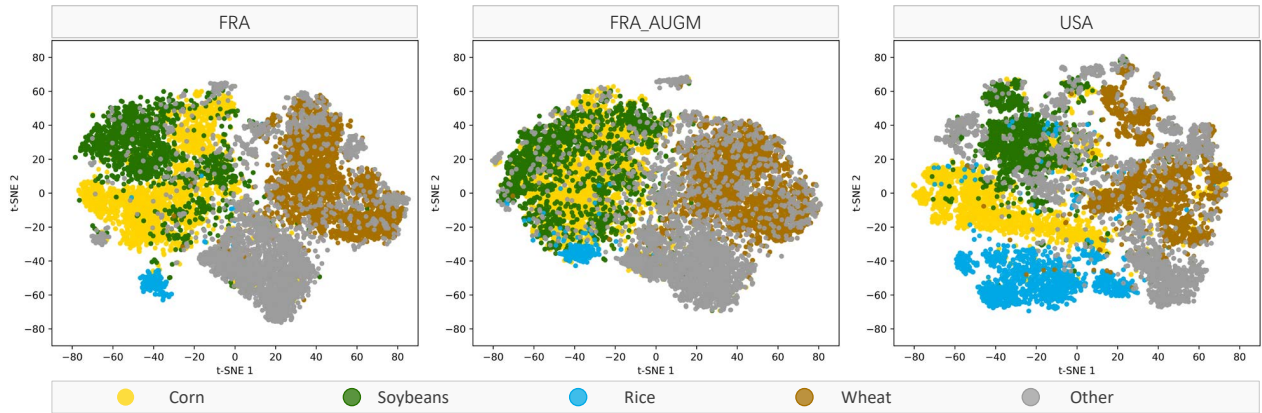


Figure 14: t-SNE visualization of S2 median features (May–November; 5-day intervals) for FRA (pre-/post-augmentation) and USA. For comparability, USA is restricted to the same five classes as FRA.

Despite the overall positive picture, we also see clear limitations. In Table 3, the results of FRA \rightarrow USA is not satisfactory, even with augmentation transformations. We visualize features via t-SNE [55] for FRA (before/after augmentation) and USA in Fig. 14. Because FRA has only five classes, we retain the same five in USA for comparability. For each group we sample 10,000 samples and visualize median features from May

to November at 5-day steps, prior to any CNN. Samples in USA look more stretched and dispersed, with classes not cleanly separated, indicating the target domain is more complex and variable. FRA, in contrast, forms tight, clean clusters: low within-class spread, but also narrow coverage of real-world variation. After augmentation, FRA_AUGM becomes wider and blurrier, with more overlap across classes. Part of this “messiness” comes from squeezing high-dimensional structure into 2D with t-SNE; part of it is real, because the augmentations increase within-class diversity. However, augmentation widens the source range but can’t recreate the full variety seen in practice. Source diversity is the ceiling. When the target combines shifts we didn’t simulate, e.g., phenology shifts, spectral offsets, and amplitude changes at the same time, the model can still fail.

Geographic coverage is another constraint. Our evaluation focuses primarily on relatively high-income countries. These settings differ systematically from many low-income regions; most notably, intercropping is uncommon in our study countries but prevalent in Sub-Saharan Africa and parts of South Asia. Yield levels and variability also differ markedly across these contexts (e.g., due to inputs, cultivars, irrigation, and field structure), which can shift spectral-temporal signatures and label noise. These differences limit truly global applicability; direct transfer to such regions will likely be more challenging.

5.2 Does Transfer Performance Correlate with Regional Yield Levels?

To examine whether crop productivity influences the invariance of spectral-temporal features across regions, we analyze the relationship between per-hectare yields and cross-region classification accuracy. The underlying hypothesis is that transferability may be higher between regions with more similar yields, particularly when both are high yielding, because in such conditions S2 reflectance is less influenced by soil, and management practices across regions tend to be more uniform.

We collect 2023 crop yield statistics from the Food and Agriculture Organization of the United Nations (<https://www.fao.org/home/en>) and summarize them in Table 8. We then compare these with classification results from selected transfer settings involving soybean, rice, and wheat, as shown in the confusion matrices in Fig. 15.

Table 8: Crop yields (kg/ha) for major crop types across countries in 2023.

Yield (kg/ha)	Corn	Soybeans	Rice	Wheat	Sugarcane	Cotton
ARG	-	1744.5	-	-	-	-
AUS	-	-	-	-	98619.5	3787.4
BEL	9828.9	-	-	8424.7	-	-
CHN	-	-	7132.5	-	-	-
FRA	9762.5	2458.4	5632.8	7201.6	-	-
GBR	-	-	-	8127.9	-	-
NLD	10265.5	-	-	8467.0	-	-
USA	11130.6	3398.7	8572.8	3269.4	79309.2	2843.0

For soybean, USA and ARG show a larger yield gap (3398.7 vs. 1744.5 kg/ha) than that between USA and FRA (3398.7 vs. 2458.4 kg/ha). Yet USA → ARG achieves higher classification accuracy than USA → FRA. This is likely due to the simplified classification setup in ARG, which only involves a binary distinction (soybean vs. other), reducing ambiguity and improving accuracy.

For rice, yield differences between USA (8572.8 kg/ha) and CHN (7132.5 kg/ha) are smaller than those between USA and FRA (5632.8 kg/ha), yet USA → FRA shows better classification accuracy than USA → CHN. This may be attributed to higher label reliability in FRA, where ground-truth labels are farmer-reported, as opposed to automatically generated labels in CHN.

For wheat, despite USA having much lower yield (3269.4 kg/ha) than FRA (7201.6 kg/ha) and NLD (8467.0 kg/ha), classification accuracy remains high in both transfer directions. This shows that, at least under the roughly twofold yield gap in this example, yield differences do not prevent high performance in cross-continent transfer.



Figure 15: Confusion matrices of 2D CropNet with data augmentation from Table 3 for four cross-region transfer settings, used to analyze the potential relationship between classification performance and crop yields.

Taken together, within the yield ranges represented in our study (mostly modest, except for approximately twofold differences in soybean and wheat), we do not find a consistent relationship between yield levels and cross-region accuracy; other factors, such as training-data quantity and quality, class balance, and the spectral-temporal distinctiveness of the crop types, appear to matter more. At the same time, our analysis cannot rule out yield as a dominant factor when gaps are much larger (e.g., in low-input smallholder systems typical of parts of Sub-Saharan Africa), where lower canopy density and greater soil/background influence may substantially alter S2 signals. Evaluating such settings requires data that span a broader yield spectrum and will be a direction for future work.

5.3 How Can Transferability Reach the Regions Most in Need?

An especially important next step is to transfer crop type mapping tools to regions that are most food insecure and in urgent need of reliable agricultural statistics. For instance, many countries in sub-Saharan Africa remain chronically data-scarce despite being highly vulnerable to food crises [56, 57]. Our results show that cross-hemisphere transfer is feasible: being in the Southern Hemisphere is not a barrier. The key question is when transfer fails.

The yield range in the CropGlobe dataset does not reach the lower levels common in many African systems: corn yields in some areas are less than 5 t/ha (CropGlobe’s lowest is France at approximately 9.8 t/ha), and rice yields are often less than 2 t/ha (France in our dataset is about 5.6 t/ha). Under such low-yield (low-biomass) conditions, combined with differences in phenology and environment, models trained on our data may not transfer directly to these more challenging regions.

Even so, our results already show that S2 2D median features achieve consistently strong cross-country accuracy within the same continent, and we hypothesize a feasible path forward. By collecting crop type data from neighboring or similar countries to supplement the training data, the model can be exposed to spectral-temporal patterns close to those expected in the target. Approaches such as soft labels [58, 59], label transfer [60], and weak constraints [61] can be used to anchor the model to the target region’s distribution and require only a small amount of local ground truth for calibration. This path does not require the same quantity or quality of labels as in our study, but can improve robustness and, while moving toward operational mapping, greatly reduce local labeling needs and costs.

The transfer learning method developed here is able to serve as a first step. On this basis, future work will extend to more challenging settings, such as lower-yield, intercropped, and highly heterogeneous systems. A longer-term aim is to enable operational crop type mapping in low-income, data-scarce countries and to lower labeling and statistical costs. This is a promising route to bring cost-effective crop type mapping to

the regions that need it most.

6 Conclusion

This study systematically investigates the geographic invariance of various remote sensing features for cross-region crop type classification. By constructing the CropGlobe dataset, which spans eight countries across five continents, we conduct experiments under three transfer settings: cross-country, cross-continent, and cross-hemisphere. The results show that the 2D median features derived from S2 exhibit significant advantages in geographic transferability. To reduce overfitting risks caused by low-dimensional inputs and enhance the geographic invariance of the learned representations, we introduce CropNet, an extremely lightweight CNN that outperforms mainstream convolutional models while maintaining the smallest parameter count. Furthermore, to mitigate spectral and phenological shifts across regions, we develop a data augmentation strategy for temporal multi-spectral features, which enhances the model’s robustness and yields accuracy improvements, especially under limited training diversity. In particular, even in the most challenging cross-hemisphere settings, classification accuracy remains above 85% across tested regions, highlighting the scalability and reliability of the proposed method. Our work provides initial evidence that global crop type mapping can be practical in many settings with fewer local ground samples, and it offers a data resource and a preliminary methodological framework for studying temporal multi-spectral feature invariance in crop type identification. Looking ahead, as the dataset expands and methods evolve, this study could help move from region-specific applications toward more broadly generalizable solutions, with the potential to support agricultural monitoring, food-security planning, and policy analysis.

References

- [1] C. Monfreda, N. Ramankutty, and J. A. Foley, “Farming the planet: 2. geographic distribution of crop areas, yields, physiological types, and net primary production in the year 2000,” *Global Biogeochemical Cycles*, vol. 22, no. 1, 2008.
- [2] I. Becker-Reshef, C. Justice, B. Barker, M. Humber, F. Rembold, R. Bonifacio, M. Zappacosta, M. Budde, T. Magadzire, C. Shitote *et al.*, “Strengthening agricultural decisions in countries at risk of food insecurity: The geoglam crop monitor for early warning,” *Remote Sensing of Environment*, vol. 237, p. 111553, 2020.
- [3] D. B. Lobell and S. M. Gourdj, “The influence of climate change on global crop productivity,” *Plant Physiology*, vol. 160, no. 4, pp. 1686–1697, 2012.
- [4] M. Haasnoot, J. H. Kwakkel, W. E. Walker, and J. Ter Maat, “Dynamic adaptive policy pathways: A method for crafting robust decisions for a deeply uncertain world,” *Global Environmental Change*, vol. 23, no. 2, pp. 485–498, 2013.
- [5] J. A. Foley, N. Ramankutty, K. A. Brauman, E. S. Cassidy, J. S. Gerber, M. Johnston, N. D. Mueller, C. O’Connell, D. K. Ray, P. C. West *et al.*, “Solutions for a cultivated planet,” *Nature*, vol. 478, no. 7369, pp. 337–342, 2011.
- [6] S. Wang, G. Azzari, and D. B. Lobell, “Crop type mapping without field-level labels: Random forest transfer and unsupervised clustering techniques,” *Remote Sensing of Environment*, vol. 222, pp. 303–317, 2019.
- [7] U.S. Department of Agriculture, National Agricultural Statistics Service, “Cropland data layer,” <https://nassgeodata.gmu.edu/CropScape/>, 2024.

- [8] Agriculture and Agri-Food Canada, “Annual crop inventory,” <https://search.open.canada.ca/openmap/ba2645d5-4458-414d-b196-6303ac06c1c9>, 2023.
- [9] UK Rural Payments Agency, “Crop map of england (crome),” <https://environment.data.gov.uk/dataset/a27312b5-d6c9-4710-ad5e-382d727c1b05>, 2023.
- [10] D. M. Kluger, S. Wang, and D. B. Lobell, “Two shifts for crop mapping: Leveraging aggregate crop statistics to improve satellite-based maps in new regions,” *Remote Sensing of Environment*, vol. 262, p. 112488, 2021.
- [11] G. Yang, X. Li, P. Liu, X. Yao, Y. Zhu, W. Cao, and T. Cheng, “Automated in-season mapping of winter wheat in china with training data generation and model transfer,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 202, pp. 422–438, 2023.
- [12] H. Kerner, S. Chaudhari, A. Ghosh, C. Robinson, A. Ahmad, E. Choi, N. Jacobs, C. Holmes, M. Mohr, R. Dodhia *et al.*, “Fields of the world: A machine learning benchmark dataset for global agricultural field boundary segmentation,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 39, no. 27, 2025, pp. 28 151–28 159.
- [13] S. Wang, F. Waldner, and D. B. Lobell, “Unlocking large-scale crop field delineation in smallholder farming systems with transfer learning and weak supervision,” *Remote Sensing*, vol. 14, no. 22, p. 5738, 2022.
- [14] J. L. Soler, T. Friedel, and S. Wang, “Combining deep learning and street view imagery to map smallholder crop types,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 20, 2024, pp. 22 202–22 212.
- [15] Y. Yan and Y. Ryu, “Exploring google street view with deep learning for crop type mapping,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 171, pp. 278–296, 2021.
- [16] H. Li, X.-P. Song, M. C. Hansen, I. Becker-Reshef, B. Adusei, J. Pickering, L. Wang, L. Wang, Z. Lin, V. Zalles *et al.*, “Development of a 10-m resolution maize and soybean map over china: Matching satellite-based crop classification with sample-based area estimation,” *Remote Sensing of Environment*, vol. 294, p. 113623, 2023.
- [17] C. Lin, L. Zhong, X.-P. Song, J. Dong, D. B. Lobell, and Z. Jin, “Early-and in-season crop type mapping without current-year ground truth: Generating labels from historical information via a topology-based approach,” *Remote Sensing of Environment*, vol. 274, p. 112994, 2022.
- [18] M. Belgiu, W. Bijker, O. Csillik, and A. Stein, “Phenology-based sample generation for supervised crop type classification,” *International Journal of Applied Earth Observation and Geoinformation*, vol. 95, p. 102264, 2021.
- [19] M. Belgiu and O. Csillik, “Sentinel-2 cropland mapping using pixel-based and object-based time-weighted dynamic time warping analysis,” *Remote Sensing of Environment*, vol. 204, pp. 509–523, 2018.
- [20] D. M. Johnson and R. Mueller, “Pre-and within-season crop type classification trained with archival land cover information,” *Remote Sensing of Environment*, vol. 264, p. 112576, 2021.
- [21] Y. Cai, K. Guan, J. Peng, S. Wang, C. Seifert, B. Wardlow, and Z. Li, “A high-performance and in-season classification system of field-level crop types using time-series landsat data and a machine learning approach,” *Remote Sensing of Environment*, vol. 210, pp. 35–47, 2018.

- [22] I. Aneece, D. Foley, P. Thenkabail, A. Oliphant, and P. Teluguntla, "New generation hyperspectral data from desis compared to high spatial resolution planetscope data for crop type classification," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 15, pp. 7846–7858, 2022.
- [23] I. Aneece and P. S. Thenkabail, "New generation hyperspectral sensors desis and prisma provide improved agricultural crop classifications," *Photogrammetric Engineering & Remote Sensing*, vol. 88, no. 11, pp. 715–729, 2022.
- [24] I. Aneece and P. Thenkabail, "Accuracies achieved in classifying five leading world crop types and their growth stages using optimal earth observing-1 hyperion hyperspectral narrowbands on google earth engine," *Remote Sensing*, vol. 10, no. 12, p. 2027, 2018.
- [25] I. Aneece and P. S. Thenkabail, "Classifying crop types using two generations of hyperspectral sensors (hyperion and desis) with machine learning on the cloud," *Remote Sensing*, vol. 13, no. 22, p. 4704, 2021.
- [26] L. Blickensdörfer, M. Schwieder, D. Pflugmacher, C. Nendel, S. Erasmi, and P. Hostert, "Mapping of crop types and crop sequences with combined time series of sentinel-1, sentinel-2 and landsat 8 data for germany," *Remote Sensing of Environment*, vol. 269, p. 112831, 2022.
- [27] Y. Zhou, J. Luo, L. Feng, Y. Yang, Y. Chen, and W. Wu, "Long-short-term-memory-based crop classification using high-resolution optical images and multi-temporal sar data," *GIScience & Remote Sensing*, vol. 56, no. 8, pp. 1170–1191, 2019.
- [28] S. Kobayashi and H. Ide, "Rice crop monitoring using sentinel-1 sar data: A case study in saku, japan," *Remote Sensing*, vol. 14, no. 14, p. 3254, 2022.
- [29] S. Di Tommaso, S. Wang, and D. B. Lobell, "Combining gedi and sentinel-2 for wall-to-wall mapping of tall and short crops," *Environmental Research Letters*, vol. 16, no. 12, p. 125002, 2021.
- [30] S. Di Tommaso, S. Wang, R. Strey, and D. B. Lobell, "Mapping sugarcane globally at 10 m resolution using gedi and sentinel-2," *Earth System Science Data Discussions*, vol. 2024, pp. 1–26, 2024.
- [31] X. Huang, A. Vrieling, Y. Dou, M. Belgiu, and A. Nelson, "A robust method for mapping soybean by phenological aligning of sentinel-2 time series," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 218, pp. 1–18, 2024.
- [32] C. Pelletier, G. I. Webb, and F. Petitjean, "Temporal convolutional neural network for the classification of satellite image time series," *Remote Sensing*, vol. 11, no. 5, p. 523, 2019.
- [33] G. Tseng, R. Cartuyvels, I. Zvonkov, M. Purohit, D. Rolnick, and H. Kerner, "Lightweight, pre-trained transformers for remote sensing timeseries," *arXiv preprint arXiv:2304.14065*, 2023.
- [34] I. Mariotto, P. Thenkabail, and I. Aneece, "Global hyperspectral imaging spectral-library of agricultural crops for central asia v001," NASA Land Processes Distributed Active Archive Center, 2020. [Online]. Available: <https://doi.org/10.5067/COMMUNITY/GHISA/GHISACASIA.001>
- [35] P. Thenkabail and I. Aneece, "Global hyperspectral imaging spectral-library of agricultural crops for conterminous united states v001," NASA Land Processes Distributed Active Archive Center, 2019. [Online]. Available: <https://doi.org/10.5067/COMMUNITY/GHISA/GHISACONUS.001>
- [36] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.

- [37] M. Tan and Q. Le, “Efficientnetv2: Smaller models and faster training,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 10 096–10 106.
- [38] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, “A convnet for the 2020s,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 11 976–11 986.
- [39] G. Tseng, I. Zvonkov, C. L. Nakalembe, and H. Kerner, “Cropharvest: A global dataset for crop-type classification,” in *35th Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021.
- [40] K. Van Tricht, J. Degerickx, S. Gilliams, D. Zanaga, M. Battude, A. Grosu, J. Brombacher, M. Lesiv, J. C. L. Bayas, S. Karanam *et al.*, “Worldcereal: a dynamic open-source system for global-scale, seasonal, and reproducible crop and irrigation mapping,” *Earth System Science Data*, vol. 15, no. 12, pp. 5491–5515, 2023.
- [41] X.-P. Song, M. C. Hansen, P. Potapov, B. Adusei, J. Pickering, M. Adami, A. Lima, V. Zalles, S. V. Stehman, C. M. Di Bella *et al.*, “Massive soybean expansion in south america since 2000 and implications for conservation,” *Nature Sustainability*, vol. 4, no. 9, pp. 784–792, 2021.
- [42] Queensland Government, “Queensland seasonal crop (qsc) mapping — summer and winter crop classifications,” <https://www.qld.gov.au/environment/land/management/mapping/statewide-monitoring/crops>, 2023.
- [43] M. J. Pringle, “Detecting the annual areal extent of sugarcane crops in queensland, australia,” *Remote Sensing Applications: Society and Environment*, vol. 22, p. 100496, 2021.
- [44] Flemish Department of Agriculture and Fisheries, “Landbouwgebruikspercelen (lgp) — annual crop parcel declarations in belgium,” <https://landbouwcijfers.vlaanderen.be/open-geodata-landbouwgebruikspercelen>, 2023.
- [45] Q. Li, H. Wang, Y. Zhang, X. Du, Y. Dong, Y. Shen *et al.*, “Mapping crop type in northeast china during 2017–2023 (including code for classification),” figshare. Dataset, 2024. [Online]. Available: <https://doi.org/10.6084/m9.figshare.25346038.v1>
- [46] Agence de Services et de Paiement, “Registre parcellaire graphique (rpg) — annual agricultural parcel database of france,” <https://geoservices.ign.fr/rpg>, 2023.
- [47] Netherlands Enterprise Agency, “Basisregistratie gewaspercelen (brp) — annual agricultural parcel and crop type register of the netherlands,” https://service.pdok.nl/rvo/brpgewaspercelen/atom/v1_0/basisregistratie-gewaspercelen-brp.xml, 2023.
- [48] J. W. Rouse, R. H. Haas, J. A. Schell, and D. W. Deering, “Monitoring vegetation systems in the great plains with erts,” in *Third Earth Resources Technology Satellite-1 Symposium*, ser. NASA Special Publication, vol. SP-351. Washington, D.C.: NASA, 1973, pp. 309–317.
- [49] T. Nguyen, M. Raghu, and S. Kornblith, “Do wide and deep networks learn the same things? uncovering how neural network representations vary with width and depth,” *arXiv preprint arXiv:2010.15327*, 2020.
- [50] D. Arpit, S. Jastrzebski, N. Ballas, D. Krueger, E. Bengio, M. S. Kanwal, T. Maharaj, A. Fischer, A. Courville, Y. Bengio *et al.*, “A closer look at memorization in deep networks,” in *International Conference on Machine Learning*. PMLR, 2017, pp. 233–242.
- [51] A. A. Gitelson, A. Viña, V. Ciganda, D. C. Rundquist, and T. J. Arkebauer, “Remote estimation of canopy chlorophyll content in crops,” *Geophysical Research Letters*, vol. 32, no. 8, 2005.

- [52] G. Forestier, F. Petitjean, H. A. Dau, G. I. Webb, and E. Keogh, “Generating synthetic time series to augment sparse datasets,” in *2017 IEEE International Conference on Data Mining*. IEEE, 2017, pp. 865–870.
- [53] J. Tompson, R. Goroshin, A. Jain, Y. LeCun, and C. Bregler, “Efficient object localization using convolutional networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 648–656.
- [54] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, “Learning deep features for discriminative localization,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2921–2929.
- [55] L. v. d. Maaten and G. Hinton, “Visualizing data using t-sne,” *Journal of Machine Learning Research*, vol. 9, no. Nov, pp. 2579–2605, 2008.
- [56] J. Sarron, D. Beillouin, J. Huat, J.-M. Koffi, J. Diatta, É. Malézieux, and E. Faye, “Digital agriculture to fulfil the shortage of horticultural data and achieve food security in sub-saharan africa,” in *IV All Africa Horticultural Congress-AAHC 2021: Transformative Innovations in Horticulture*, 2021, pp. 239–246.
- [57] D. Lee, W. Anderson, X. Chen, F. Davenport, S. Shukla, R. Sahajpal, M. Budde, J. Rowland, J. Verdin, L. You *et al.*, “Harveststat africa—harmonized subnational crop statistics for sub-saharan africa,” *Scientific Data*, vol. 12, no. 1, p. 690, 2025.
- [58] S. Wang, W. Chen, S. M. Xie, G. Azzari, and D. B. Lobell, “Weakly supervised deep learning for segmentation of remote sensing imagery,” *Remote Sensing*, vol. 12, no. 2, p. 207, 2020.
- [59] X.-Y. Tong, R. Dong, and X. X. Zhu, “Global high categorical resolution land cover mapping via weak supervision,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 220, pp. 535–549, 2025.
- [60] X. Zhu, Z. Ghahramani, and J. D. Lafferty, “Semi-supervised learning using gaussian fields and harmonic functions,” in *Proceedings of the 20th International Conference on Machine Learning*, 2003, pp. 912–919.
- [61] K. Ganchev, J. Graça, J. Gillenwater, and B. Taskar, “Posterior regularization for structured latent variable models,” *The Journal of Machine Learning Research*, vol. 11, pp. 2001–2049, 2010.