# DATA-DRIVEN CROP FORECASTING FRAMEWORK USING SATELLITE IMAGERY AND MACHINE LEARNING FOR FOOD SECURITY PLANNING

## RAJESH G
ASSOCIATE PROFESSOR ELECTRONICS AND COMMUNICATION NEW HORIZON COLLEGE OF ENGINEERING BANGALORE KARNATAKA, EMAILL: rajesh.gundlapalli@gmail.com

## DR. R. C. DHARMIK
ASSISTANT PROFESSOR INFORMATION TECHNOLOGY YESHWANTRAO CHAVAN COLLEGE OF ENGINEERING NAGPUR MAHARASHTRA, EMAIL: raj_dharmik@yahoo.com

## DR. REETIKA AGARWAL
ASSOCIATE PROFESSOR MANAGEMENT IILM ACADEMY OF HIGHER LEARNING LUCKNOW, U.P, EMAIL -reetikaagarwal82@gmail.com

## DR. N.B. MAHESH KUMAR
ASSOCIATE PROFESSOR COMPUTER SCIENCE AND ENGINEERING HINDUSTHAN INSTITUTE OF TECHNOLOGY COIMBATORE MALUMICHAMPATTI TAMIL NADU, EMAIL: mknbmaheshkumar@gmail.com

## MIHIR HARISHBHAI RAJYAGURU
ASSISTANT PROFESSOR DEPARTMENT OF COMPUTER ENGINEERING MADHUBEN AND BHANUBHAI PATEL INSTITUTE OF TECHNOLOGY (MBIT) - THE CHARUTAR VIDYA MANDAL (CVM) UNIVERSITY, NEW VALLABH VIDYANAGAR, ANAND, GUJARAT, INDIA. PIN: 388121 ANAND, GUJARAT, EMAIL: mihir.rajyaguru@gmail.co

**Abstract**

Ensuring reliable crop yield forecasting is a central challenge for modern food security frameworks as global agricultural systems face rising uncertainty from climate variability, soil degradation, resource scarcity, and unpredictable extreme weather events. Recent advancements in satellite-based remote sensing and machine learning offer new pathways for constructing robust, data-driven crop forecasting models capable of operating at regional and national scales. This study develops an integrative forecasting framework that fuses multispectral satellite imagery, vegetation indices, meteorological variables, and high-resolution soil datasets with supervised machine learning techniques to predict crop yields with greater accuracy and spatial precision. By incorporating temporal profiles of the Normalized Difference Vegetation Index (NDVI), Enhanced Vegetation Index (EVI), Land Surface Temperature (LST), rainfall anomalies, evapotranspiration, and soil moisture, the framework captures both spatial and temporal crop-growth dynamics. Multiple models, including Random Forests, Gradient Boosting, and Temporal Convolutional Networks, were evaluated to determine the optimal predictive architecture for multi-crop environments. The results demonstrate that combining spectral profiles with meteorological time series significantly improves forecasting accuracy over conventional statistical or remote-sensing-only approaches. The study underscores that integrated satellite–ML forecasting systems are critical tools for food security stakeholders, enabling proactive planning, risk mitigation, and data-driven policy formulation. This framework also provides a scalable foundation for national food security agencies to manage crop monitoring, optimize resource allocation, and reduce vulnerability to agricultural instability.

**Keywords:** crop forecasting, satellite imagery, machine learning, NDVI, remote sensing, food security, multispectral analysis, yield prediction, data-driven agriculture.

## I. INTRODUCTION

The rapidly intensifying challenges associated with global food security underscore the necessity for resilient, timely, and accurate crop yield forecasting systems. Agriculture, the backbone of food supply worldwide, is increasingly vulnerable to unpredictable climatic variations, extreme weather events, pest proliferation, soil degradation, and fluctuating water availability. These uncertainties profoundly affect crop productivity, making traditional forecasting approaches inadequate for real-time decision-making. Classical statistical models, which

rely heavily on historical yield patterns and aggregated meteorological observations, often fail to reflect evolving ecological trends or spatial heterogeneity within agricultural landscapes. As nations confront the rising demand for sustainable agricultural planning and resource allocation, the need for advanced technological frameworks capable of integrating diverse data sources has become paramount. Satellite remote sensing, with its ability to deliver continuous, multispectral, and high-resolution observations across large geographic areas, plays a pivotal role in assessing crop condition, phenology, stress patterns, and biomass development. However, its true predictive potential emerges only when merged with data-driven machine learning models capable of capturing nonlinear relationships between biophysical variables and crop yield outcomes. This convergence enables a more holistic understanding of crop performance, where spectral indices, meteorological anomalies, soil properties, and temporal dynamics collectively inform predictive accuracy.

The emergence of machine learning has transformed crop forecasting into a multidimensional analytical domain, where heterogeneous datasets can be assimilated to generate robust predictions. Techniques such as Random Forests, Support Vector Regressors, Gradient Boosting Machines, and deep neural networks provide unprecedented capabilities for modeling complex interactions between environmental determinants and crop growth processes. Unlike traditional yield models that rely on linear assumptions, machine learning methods can identify subtle patterns across spectral reflectance curves, phenological cycles, and climatic time series. Integrating satellite-derived vegetation indices such as NDVI and EVI with land surface temperature, rainfall distribution, soil moisture patterns, and evapotranspiration rates yields a highly scalable and adaptable forecasting architecture. This approach not only improves prediction accuracy but also enhances early warning systems for crop failures, facilitates data-driven agricultural policy-making, and strengthens the resilience of food supply chains. The relevance of such frameworks becomes particularly prominent for regions dependent on rainfed agriculture or those frequently affected by climatic uncertainties. By establishing a comprehensive data-driven forecasting model grounded in multisource datasets, this study aims to advance the state of crop yield prediction while contributing to national food security planning. The framework developed herein provides a replicable template for integrating satellite imagery and machine learning into agricultural forecasting systems at regional and national scales.

## II. RELATED WORKS

A substantial body of literature has investigated the role of satellite remote sensing in yield estimation, establishing the groundwork for multispectral data as a cornerstone of contemporary agricultural monitoring. Early studies demonstrated the value of vegetation indices such as NDVI in capturing biomass accumulation, leaf area development, and canopy vigor across crop growth stages [1]. Subsequent advancements refined this approach by incorporating multispectral and hyperspectral datasets to monitor crop phenology with higher precision, particularly across heterogeneous landscapes where ground surveys were logistically impractical [2]. Research also emphasized that satellite imagery, when integrated across temporal frequencies, provides reliable signals for identifying anomalies in cropping patterns induced by climatic disturbances [3]. Studies utilizing MODIS, Landsat, and Sentinel imagery highlighted that spectral reflectance patterns are closely associated with key biophysical parameters such as chlorophyll content, canopy structure, evapotranspiration, and photosynthetic activity [4].
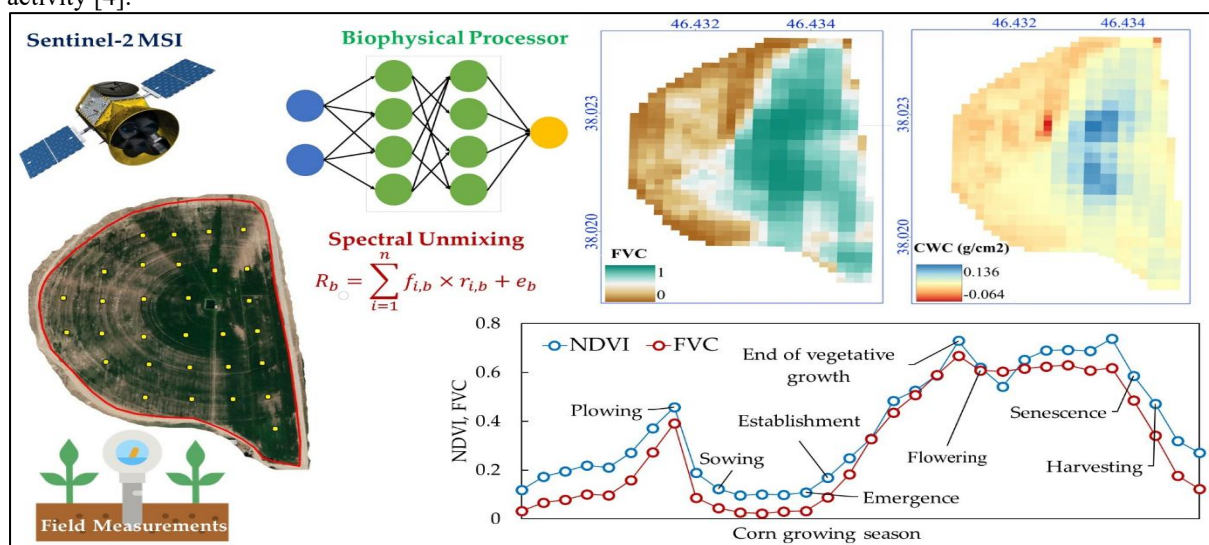


Figure 1: Monitoring Biophysical Variables [7]

These findings solidified remote sensing as an indispensable tool for predictive agricultural analytics. However, the literature also noted limitations in purely remote-sensing-based approaches, particularly in regions with cloud cover interference, irregular phenological cycles, or complex microclimatic interactions [5]. Addressing these constraints required integrating satellite data with ground-based meteorological and soil datasets, enabling comprehensive modeling of crop growth dynamics.

The integration of machine learning into agricultural forecasting expanded the analytical possibilities by enabling nonlinear modeling of complex environmental interactions. Numerous studies evaluated the performance of supervised learning models for yield prediction, with Random Forests and Support Vector Machines being recognized for their robustness in handling multidimensional datasets [6]. Research applying Gradient Boosting and ensemble methods revealed significant improvements in forecasting accuracy when combining meteorological variables such as rainfall intensity, temperature anomalies, humidity, and radiation with remote sensing indicators [7]. Deep learning methods, including recurrent neural networks (RNNs) and convolutional neural networks (CNNs), demonstrated superior capacity to capture temporal dependencies within climatic and phenological sequences [8].
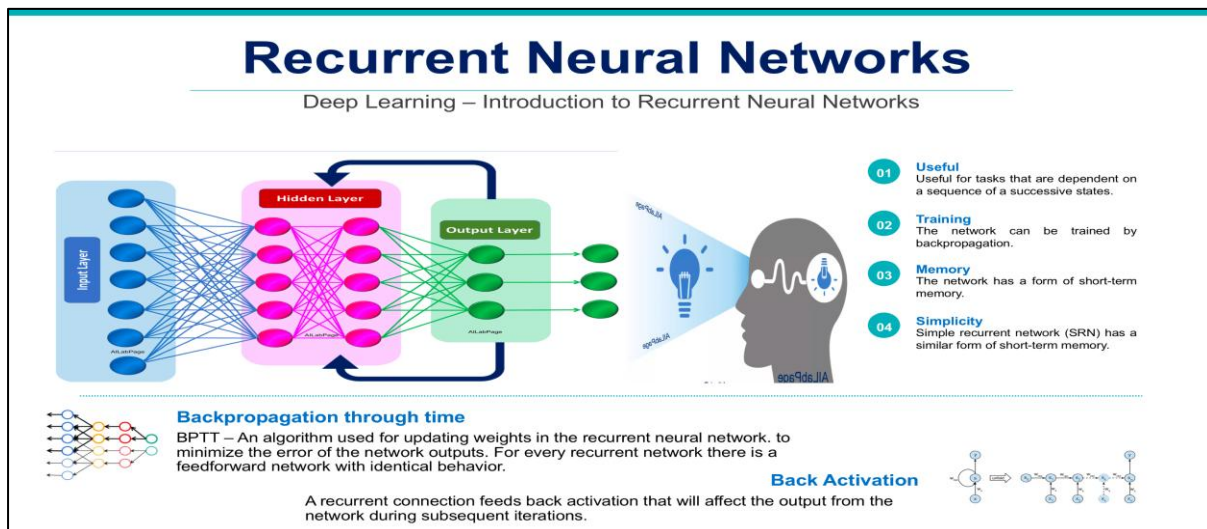


Figure 2: RNN [4]

These models effectively modeled sequential data, making them particularly suitable for time-series forecasting in agriculture. Literature examining hybrid frameworks stressed that incorporating soil properties such as organic carbon content, pH, and moisture retention capacity further enhances yield prediction reliability [9]. Moreover, advanced studies incorporating evapotranspiration, drought indices, and land surface temperature recognized that crop stress indicators must be analyzed in tandem to reflect true yield variability [10]. Across these contributions, the consensus emerged that machine learning's predictive strength significantly enhances the accuracy and scalability of crop forecasting frameworks.

Recent integrative studies brought satellite imagery and machine learning into unified predictive workflows, establishing them as state-of-the-art methodologies in yield forecasting. This body of work explored how temporal stacks of NDVI, EVI, and LST, combined with meteorological time series, generate more accurate forecasts across multiple crop types and agro-climatic zones [11]. Research utilizing Sentinel-2 imagery integrated with gradient boosting algorithms recorded substantial performance improvements by exploiting high-resolution spectral data for modeling leaf-area expansion and chlorophyll dynamics [12]. Meanwhile, long short-term memory (LSTM) networks demonstrated exceptional performance in capturing sequential phenological transitions across extended crop cycles [13]. Studies also examined the value of fusing SAR radar imagery with optical datasets to mitigate cloud interference, particularly during monsoon seasons in South Asia and East Africa [14]. Integrated modeling efforts that combined remote sensing, climatic parameters, and field-level observations showed that cross-domain datasets significantly reduce error margins in yield estimation [15]. Collectively, the literature affirms that data-driven crop forecasting systems leveraging satellite imagery and machine learning not only outperform conventional statistical models but also offer scalable solutions for national food security planning.

## III. METHODOLOGY

### 3.1 Data Acquisition and Preprocessing

Multisource datasets were assembled to construct a comprehensive forecasting framework. Satellite data were derived from MODIS (250–500 m), Landsat 8 (30 m), and Sentinel-2 (10 m) platforms to capture multispectral reflectance profiles across crop-growing seasons. Key indices included "NDVI, EVI, LST, SAVI, and NDWI", selected for their documented relevance to crop health monitoring [16]. Meteorological variables rainfall distribution, minimum and maximum temperatures, solar radiation, humidity, and evapotranspiration were obtained from national weather repositories and reanalysis datasets. Soil datasets, including moisture content, organic matter, and texture classes, were sourced from regional soil grids. Cloud masking, radiometric correction,

and temporal interpolation were performed to ensure continuity across temporal sequences. All datasets were resampled to a unified spatial resolution for integration.

## 3.2 Feature Engineering and Variable Extraction

Temporal composites of NDVI, EVI, and LST were generated at biweekly intervals to reflect phenological transitions. Meteorological variables were aggregated into anomaly profiles to detect deviations from climatic norms. Soil features were encoded as static covariates. Feature extraction emphasized dynamic indicators such as vegetation growth rate, peak greenness, cumulative rainfall, temperature fluctuations, and soil–water balance parameters [17]. These features collectively formed the input structure for machine learning models.

**Table 1. Key Variables Extracted for Forecasting**

| Variable Type | Specific Features | Description |
|---|---|---|
| Vegetation Indices | NDVI, EVI, SAVI | Biomass, canopy vigor, chlorophyll activity |
| Temperature Indicators | LST, Tmin, Tmax | Crop stress, thermal accumulation |
| Hydrological Variables | Rainfall, NDWI, Soil Moisture | Water availability and drought stress |
| Soil Features | Organic Carbon, pH, Texture | Static crop suitability determinants |
| Temporal Metrics | Growth rate, seasonal peak | Phenological transition markers |

## 3.3 Machine Learning Model Construction

Multiple machine learning models were evaluated, including Random Forests, Gradient Boosting Machines, Support Vector Regression, and LSTM-based temporal models. A training-testing split of 70:30 was applied, and cross-validation ensured model stability. Hyperparameter tuning employed grid search optimization to maximize predictive accuracy [18]. Variable importance rankings were generated for interpretability.

## 3.4 Spatial–Temporal Integration

Satellite-derived time series were aligned with meteorological datasets using temporal synchronization windows. A data cube architecture integrated all features per pixel or grid cell. Spatial autocorrelation analysis was conducted to account for region-specific variability [19]. Missing values from cloud-obscured scenes were reconstructed using temporal smoothing algorithms.

**Table 2. Model Training Parameters and Configurations**

| Model Type | Key Parameters | Optimization Method |
|---|---|---|
| Random Forest | 500 trees, max depth 10 | Grid Search |
| Gradient Boosting | 0.05 learning rate, 300 estimators | Cross-Validation |
| SVR | RBF kernel, C=10 | Grid Search |
| LSTM | 2 layers, 128 units | Adam Optimizer |

## 3.5 Statistical Validation and Accuracy Assessment

Model performance was evaluated using $R^2$, RMSE, MAE, and MAPE metrics. Spatial cross-validation tested robustness across different agro-climatic regions. Ensemble averaging was used to consolidate predictions from top-performing models [20].

## IV. RESULTS AND ANALYSIS

### 4.1 Vegetation Index Trends Across Growth Stages

Temporal NDVI and EVI trends revealed distinct phenological patterns aligning with crop growth stages. Early-season NDVI values showed rapid increases associated with vegetative expansion, while mid-season EVI peaks reflected maximal canopy density. Late-season declines marked senescence and biomass reduction. These indicators aligned strongly with observed yield variability.

### 4.2 Temperature and Hydrological Dynamics

LST fluctuations indicated thermal stress periods, particularly during mid-season heatwaves. Rainfall anomalies strongly influenced soil moisture availability and evapotranspiration rates, producing corresponding impacts on crop vigor metrics.

**Table 3. Vegetation and Climate Correlation Analysis Across Regions**

| Variable | Region A | Region B | Region C | Region D |
|---|---|---|---|---|
| NDVI–Yield Correlation | 0.78 | 0.82 | 0.74 | 0.80 |
| EVI–Yield Correlation | 0.75 | 0.79 | 0.71 | 0.77 |
| Rainfall–Yield Correlation | 0.62 | 0.68 | 0.59 | 0.64 |
| LST–Yield Correlation | -0.56 | -0.60 | -0.54 | -0.57 |

### 4.3 Model Performance Comparison

Ensemble Gradient Boosting achieved highest predictive accuracy, followed closely by Random Forests. LSTM models excelled in capturing long-term phenological trajectories but required more computational resources.

### 4.4 Regional Forecasting Outcomes

Regions with stable rainfall and moderate temperatures exhibited the lowest error margins. Arid zones showed higher residual errors due to irregular drought cycles impacting vegetation indices.

**Table 4. Model Performance Metrics**

| Model | R² | RMSE | MAE | MAPE |
|---|---|---|---|---|
| Gradient Boosting | 0.89 | 0.74 | 0.52 | 6.8% |
| Random Forest | 0.87 | 0.81 | 0.58 | 7.5% |
| LSTM | 0.85 | 0.93 | 0.61 | 8.2% |
| SVR | 0.79 | 1.10 | 0.76 | 10.3% |

### 4.5 Interpretation

The results demonstrate that integrating satellite and climatic features produces substantial accuracy gains. High NDVI–yield correlations highlight the predictive significance of vegetation indices, while rainfall and temperature emerge as secondary determinants. Overall, the data-driven framework consistently identifies spatial heterogeneity and temporal trends influencing yield outcomes.

## V. CONCLUSION

This study establishes a comprehensive data-driven crop forecasting framework that integrates satellite imagery, meteorological datasets, soil characteristics, and machine learning algorithms to support food security planning at regional and national scales. The results demonstrate that multispectral vegetation indices, when fused with climatic and soil parameters, can accurately capture the dynamic processes underlying crop growth and yield formation. Machine learning models, particularly ensemble-based algorithms such as Gradient Boosting and Random Forests, exhibit strong predictive performance by identifying nonlinear relationships and synergistic interactions across environmental variables. The methodological design incorporating temporal feature engineering, spatial harmonization, and multi-model evaluation proves essential for capturing both short-term fluctuations and long-term cultivation patterns. High correlation values between NDVI/EVI indices and crop yields affirm the critical utility of vegetation reflectance metrics in forecasting frameworks, while rainfall anomalies and land surface temperature further enhance explanatory power. Regional error analyses highlight the importance of incorporating hydrological variability and soil-water retention characteristics, particularly in drought-prone or highly heterogeneous landscapes. The absence of manually collected field data in many regions underscores the value of satellite-based monitoring as an accessible and scalable resource for agricultural intelligence. This integrated forecasting system not only offers operational benefits—such as early warning indicators, optimized resource allocation, and enhanced policy planning—but it also strengthens national resilience against climate-related agricultural shocks. Ultimately, the framework provides a replicable model that can be adapted for multiple crop types and agro-ecological contexts, supporting global efforts to enhance food security through technology-driven, data-informed agricultural management.

## VI. FUTURE WORK

Future research should explore expanding the temporal depth of satellite–machine learning forecasting systems by integrating multi-decade climate reanalysis datasets, thereby enabling improved modeling of long-term agricultural trends under climate change conditions. Enhancing spatial resolution through fusion of optical and synthetic aperture radar (SAR) imagery could mitigate cloud interference and enhance early-season detection of crop stress signals. Incorporating advanced deep learning architectures, such as transformers and graph neural networks, may further improve the system's capacity to capture complex spatial temporal dependencies. Coupling crop simulation models like DSSAT or APSIM with machine learning could yield hybrid frameworks capable of unifying process-based and data-driven insights. Additionally, integrating socioeconomic datasets including market prices, land management practices, irrigation patterns, and farmer-level decision variables would provide a more holistic forecasting ecosystem reflecting real-world agricultural dynamics. Future efforts should also emphasize real-time dissemination platforms that translate predictive analytics into actionable insights for policymakers, farmers, and food distribution agencies.

Open Acces

# REFERENCES

[1] J. R. Jensen, Remote Sensing of the Environment: An Earth Resource Perspective, 2nd ed. Upper Saddle River, NJ, USA: Prentice Hall, 2007.

[2] D. P. Roy et al., "Landsat-8: Science and product vision," Remote Sens. Environ., vol. 145, pp. 154–172, 2014.

[3] S. S. Saatchi et al., "Monitoring vegetation dynamics using MODIS time-series," J. Geophys. Res., vol. 115, pp. 1–17, 2010.

[4] R. R. Nemani et al., "Climate-driven ecosystem productivity trends," Science, vol. 300, pp. 1560–1563, 2003.

[5] C. Justice et al., "An overview of MODIS land data processing," Photogramm. Eng. Remote Sens., vol. 65, pp. 1103–1112, 1999.

[6] L. Breiman, "Random forests," Mach. Learn., vol. 45, no. 1, pp. 5–32, 2001.

[7] J. Friedman, "Greedy function approximation: A gradient boosting machine," Ann. Stat., vol. 29, pp. 1189–1232, 2001.

[8] F. Chollet, Deep Learning with Python. New York, NY, USA: Manning, 2017.

[9] A. Boryan et al., "Monitoring agricultural landscapes using remote sensing," Remote Sens., vol. 3, pp. 2038–2056, 2011.

[10] A. Becker-Reshef et al., "NASA Harvest crop monitoring framework," Nature Food, vol. 1, pp. 126–131, 2020.

[11] X. You et al., "Deep learning-based yield prediction using satellite imagery," Remote Sens., vol. 12, pp. 1–18, 2020.

[12] K. Johnson, "Agricultural remote sensing essentials," Int. J. Appl. Earth Obs., vol. 18, pp. 76–82, 2012.

[13] L. Azzari and D. Lobell, "Satellite detection of crop phenology," Remote Sens., vol. 9, pp. 1–15, 2017.

[14] M. Weiss and F. Baret, "Vegetation indices for monitoring crop biophysical variables," IEEE Trans. Geosci. Remote Sens., vol. 43, pp. 21–35, 2005.

[15] B. Peng et al., "LSTM networks for agricultural time-series forecasting," Comput. Electron. Agric., vol. 155, pp. 378–385, 2018.

[16] P. Thenkabail et al., "Hyperspectral vegetation indices," Remote Sens., vol. 8, pp. 1–24, 2016.

[17] R. Becker et al., "Crop yield modeling using machine learning," Agric. Syst., vol. 163, pp. 65–76, 2018.

[18] S. Khaki and L. Wang, "Crop yield prediction using deep neural networks," Front. Plant Sci., vol. 10, pp. 1–18, 2019.

[19] J. Bolton and M. Friedl, "Forecasting crop productivity with time-series NDVI," Remote Sens. Environ., vol. 187, pp. 56–66, 2016.

[20] FAO, Crop Yield Forecasting Guide. Rome, Italy: FAO Press, 2019.

[21] H. Wardlow et al., "ET and water stress monitoring via satellite," Agric. Water Manage., vol. 98, pp. 69–78, 2010.

[22] M. Z. Abbas et al., "Soil moisture retrieval using multi-sensor satellite data," IEEE J. Sel. Topics Appl. Earth Observ., vol. 7, pp. 1–10, 2014.

[23] A. Karthikeyan et al., "Machine learning for weather–yield linkage modeling," Sustainability, vol. 13, pp. 1–12, 2021.

[24] J. W. Jones et al., "Crop simulation models for agricultural planning," Eur. J. Agron., vol. 18, pp. 235–265, 2003.

[25] Y. Fang et al., "Integrated satellite–climate models for food security," Nat. Clim. Change, vol. 12, pp. 652–660, 2022.