# Annual 30 m soybean yield mapping in Brazil using long-term satellite observations, climate data and machine learning

Xiao-Peng Song[1], Haijun Li[1], Peter Potapov[2], and Matthew C Hansen[2]

[1]Texas Tech University
[2]University of Maryland

November 22, 2022

## Abstract

Long-term spatially explicit information on crop yield is essential for understanding food security in a changing climate. Here we present a study that combines twenty-years of Landsat and MODIS data, climate and weather records, municipality-level crop yield statistics, random forests and linear regression models for mapping crop yield in a multi-temporal, multi-scale modeling framework. The study was conducted for soybean in Brazil, the world's largest producer and exporter of this commodity crop. Using a recently developed 30 m resolution, annual (2001-2019) soybean classification map product, we aggregated multi-temporal phenological metrics derived from Landsat and MODIS data over soybean pixels to the municipality scale. We combined phenological metrics with topographic features, long-term climate data, in-season weather data and soil variables as inputs to machine learning models. We trained a multi-year random forests model using yield statistics as reference and subsequently applied linear regression to adjust the biases in the direct output of the random forests model. This model combination achieved the best performance with a root-mean-square-error (RMSE) of 344 kg/ha (12% relative to long-term mean yield) and an r2 of 0.69, on the basis of 20% withheld test data. The RMSE of the leave-one-year-out assessment ranged from 259 kg/ha to 816 kg/ha. To eliminate the artifacts caused by the coarse-resolution climate and weather data, we developed multiple models with different categories of input variables. Employing the per-pixel uncertainty estimates of different models, the final soybean yield maps were produced through per-pixel model composition. We applied the models trained on 2001-2019 data to 2020 data and produced a soybean yield map for 2020, demonstrating the predictive capability of trained machine learning models for operational yield mapping in future years. Our research showed that combining satellite, climate and weather data and machine learning could effectively map crop yield at high resolution, providing critical information to understand yield growth, anomaly and food security.

1 **Annual 30 m soybean yield mapping in Brazil using long-term satellite**

2 **observations, climate data and machine learning**

3

4 **Xiao-Peng Song[1,*], Haijun Li[1] , Peter Potapov[2], Matthew C. Hansen[2]**

5     1. Department of Geosciences, Texas Tech University, Lubbock, TX, USA

6     2. Department of Geographical Sciences, University of Maryland, College Park, MD, USA

7 *Correspondence to: xiaopeng.song@ttu.edu

8

9 **Abstract**

10 Long-term spatially explicit information on crop yield is essential for understanding food security in a

11 changing climate. Here we present a study that combines twenty-years of Landsat and MODIS data,

12 climate and weather records, municipality-level crop yield statistics, random forests and linear regression

13 models for mapping crop yield in a multi-temporal, multi-scale modeling framework. The study was

14 conducted for soybean in Brazil, the world's largest producer and exporter of this commodity crop. Using

15 a recently developed 30 m resolution, annual (2001-2019) soybean classification map product, we

16 aggregated multi-temporal phenological metrics derived from Landsat and MODIS data over soybean

17 pixels to the municipality scale. We combined phenological metrics with topographic features, long-term

18 climate data, in-season weather data and soil variables as inputs to machine learning models. We trained a

19 multi-year random forests model using yield statistics as reference and subsequently applied linear

20 regression to adjust the biases in the direct output of the random forests model. This model combination

21 achieved the best performance with a root-mean-square-error (RMSE) of 344 kg/ha (12% relative to long-

22 term mean yield) and an $r^2$ of 0.69, on the basis of 20% withheld test data. The RMSE of the leave-one-

23   year-out assessment ranged from 259 kg/ha to 816 kg/ha. To eliminate the artifacts caused by the coarse-

24   resolution climate and weather data, we developed multiple models with different categories of input

25   variables. Employing the per-pixel uncertainty estimates of different models, the final soybean yield maps

26   were produced through per-pixel model composition. We applied the models trained on 2001-2019 data

27   to 2020 data and produced a soybean yield map for 2020, demonstrating the predictive capability of

28   trained machine learning models for operational yield mapping in future years. Our research showed that

29   combining satellite, climate and weather data and machine learning could effectively map crop yield at

30   high resolution, providing critical information to understand yield growth, anomaly and food security.

31   **Keywords**

32   Crop yield map; Random forests; Landsat; MODIS; Climate; Weather

33

## 1. Introduction

Reliable and timely information on crop production can inform commodity markets, insurance companies, and policy interventions in response to natural disasters and human conflict (Benami et al. 2021; Li et al. 2022; Vroege et al. 2021). Estimating crop production over a spatial unit requires information on crop harvested area and crop yield (i.e. production per unit area). Both harvested area and yield can be derived from statistical field surveys or from satellite observations (Mulla 2013; Weiss et al. 2020) . While many methods exist in mapping crop type and estimating crop area using remote sensing (e.g. Defourny et al. 2019; Gallego 2004; Gonzáles-Alonso and Cuevas 1993; Hu et al. 2021; King et al. 2017; Massey et al. 2017; Skakun et al. 2017; Song et al. 2017; Wardlow and Egbert 2008), studies are increasingly investigating direct mapping of crop yield using remote sensing data. Crop yield maps can facilitate a number of research or practical applications, such as climate impact evaluation and yield gap analysis (Lobell 2013).

Mapping crop yield requires crop type masks as a prerequisite. When crop type masks are available, two different strategies are commonly used to produce spatially explicit information on yield: the model-data integration approach and the remote sensing-based empirical approach. The model-data integration approach seeks to integrate crop simulation models with remote-sensing-derived biophysical variables for yield forecasting (Delécolle et al. 1992; Moulin et al. 1998). Crop simulation models are developed using comprehensive measurements recorded at the plot or field level, such as crop cultivar, sowing date, soil property, water and nutrient inputs, weather, and plant physiological and morphological features (e.g. leaf area index or LAI) (de Wit et al. 2019; Holzworth et al. 2014; Jones et al. 2003; Williams et al. 1989; Yang et al. 2004). The modeled processes of crop growth can be used to predict crop productivity and to evaluate the impacts of agricultural management and environmental stressors. Various techniques have been proposed to "spatialize" crop process models using time-series of satellite-based soil, plant and environmental variables, such as soil moisture, normalized difference vegetation index (NDVI), LAI, green area index (GAI), and fraction of photosynthetically active radiation (fPAR) (Battude et al. 2016;

59  Claverie et al. 2012; de Wit et al. 2012; Doraiswamy et al. 2004; Duchemin et al. 2008; Huang et al.

60  2015; Ines et al. 2013; Kang and Özdoğan 2019; Nearing et al. 2012). Yet, a general limitation of

61  applying crop process models over large areas is the lack of sufficient and accurate information about

62  model inputs (Duchemin et al. 2008; Jin et al. 2018). Moreover, the model-data integration approach

63  usually does not serve the purpose of high-resolution yield mapping. The computational cost of per-pixel

64  crop simulation is high, but such barriers are being lifted by the recent development of cloud-computing

65  platforms such as Google Earth Engine (Gorelick et al. 2017).

66  The remote sensing-based empirical approach for crop yield mapping employs regression or machine

67  learning techniques to relate vegetation variables at key crop growth stages directly to yield. An early

68  work by Tucker et al. (1980) showed that time-integrated NDVI had significant correlation with grain

69  yield in a winter wheat field in Beltsville, Maryland. Becker-Reshef et al. (2010) demonstrated that

70  seasonal peak NDVI from the Moderate Resolution Imaging Spectroradiometer (MODIS) strongly

71  correlated with winter wheat yield in Kansas and Ukraine. Franch et al. (2015) extended the Becker-

72  Reshef et al. (2010) approach by including Growing Degree Day (GDD) information, which enabled yield

73  forecasting at about one month prior to peak NDVI. Funk and Budde (2009) found that time-integrated

74  MODIS NDVI adjusted to the onset of the rainy season correlated well with maize production in

75  Zimbabwe. Yield estimation may be improved by incorporating explicit phenology information using

76  other vegetation indices beyond NDVI. Building on the work of Funk and Budde (2009), Bolton and

77  Friedl (2013) suggested that MODIS-based two-band Enhanced Vegetation Index (EVI2) standardized by

78  the greenup date correlated better than NDVI with county-level yield for maize, but indifferent for

79  soybean, over central US. Similarly, Sakamoto et al. (2013) applied a phenology detection method to

80  identify corn silking stage and demonstrated that MODIS-derived Wide Dynamic Range Vegetation

81  Index (WDRVI) (Gitelson 2004) at that stage had high correlations with yield over major corn producing

82  states of the US. Johnson (2014) proved that daytime land surface temperature (LST) negatively

83  correlated with maize and soybean yield in the US while MODIS peak NDVI positively correlated with

84    yield. Recently, Skakun et al. (2021) investigated the utility of Landsat-8, Sentinel-2, WorldView-3 and

85    Planet data for corn and soybean yield mapping over a number of sample sites in Iowa, and found that

86    surface reflectance from red-edge bands performed better than vegetation indices to reveal field-level

87    yield variability. Lobell et al. (2015) developed an approach that used simulations from a crop model to

88    train a regression to predict yields from satellite observations, and the approach was tested in industrial as

89    well as smallholder systems (Jin et al. 2019).

90    While regression-based methods are straightforward to implement, more complex algorithms and data

91    analytic techniques such as machine learning algorithms are being increasingly investigated. Using NDVI

92    from the Advanced Very High Resolution Radiometer (AVHRR) and MODIS, Li et al. (2007) compared

93    multivariate linear regression and artificial neural networks for modeling corn and soy yield over a

94    number of sample counties in the US corn belt. Likewise, Johnson et al. (2016) compared the

95    performance of multiple linear regression and nonlinear Bayesian neural networks and model-based

96    recursive partitioning for forecasting barley, canola and spring wheat yields on the Canadian Prairies.

97    Based on the finding that NDVI and LST highly correlated with crop yield, Johnson (2014) built a

98    regression tree model using multiple years of county-level yield statistics as reference and applied the

99    model to MODIS data to forecast corn and soybean yield at 250 m resolution in the US. Cai et al. (2019)

100    tested the utility of the enhanced vegetation index (EVI) from MODIS and solar-induced chlorophyll

101    fluorescence from GOME-2 and SCIAMACHY, and regression and machine learning algorithms for

102    wheat yield prediction in Australia, and found that the combination of MODIS EVI, climate data and

103    support vector machines (SVM) could achieve high performance in yield prediction. Mateo-Sanchis et al.

104    (2019) proposed a multi-sensor metric, namely the time lag between MODIS EVI and vegetation optical

105    depth (VOD) from the Soil Moisture Active Passive (SMAP) satellite, as input to nonlinear kernel ridge

106    regression for modeling county-scale crop yield in the US corn belt. Deep learning algorithms are also

107    being explored in yield estimation. Schwalbert et al. (2020) developed a method for in-season soybean

108    yield forecasting using the Long-Short Term Memory (LSTM) algorithm, MODIS-based NDVI, EVI and
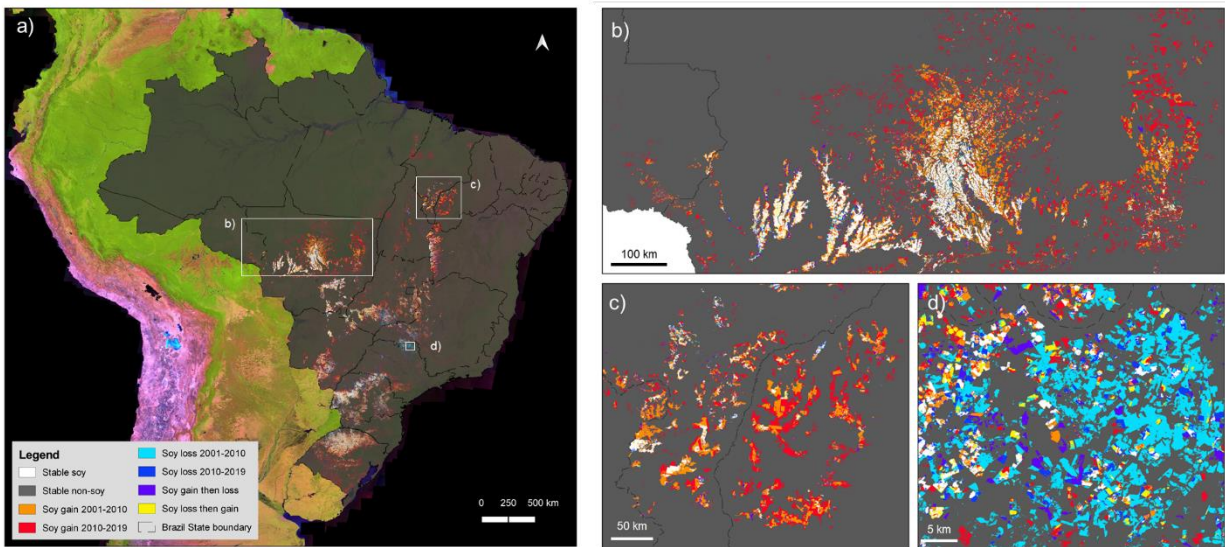
109    LST data, and precipitation data at the municipality scale in the Brazilian state of Rio Grande do Sul.

110    Recent research has also started to combine machine learning and crop models by incorporating output

111    variables from crop models as input features to machine learning algorithms for yield estimation (Paudel

112    et al. 2021; Shahhosseini et al. 2021).

113    These previous studies clearly show that crop yield estimation represents a continually active line of

114    research in remote sensing. The primary goal is to improve the accuracy of yield estimation using new

115    data and techniques, and/or to advance the date of in-season forecasting. However, most previous studies

116    are demonstrative research with limited spatial extents and/or temporal span in their study areas. Studies

117    exploring the long-term satellite data archives to evaluate the variability of crop yields also exist albeit

118    over small study areas (e.g. Gao et al. 2018; Liu et al. 2020). More importantly, common to most yield

119    mapping studies, crops in the temperate climate zone are often the target crops and target regions. Long-

120    term, large-area crop yield mapping in the tropics does not exist. Unlike the temperate region where

121    climate conditions are relatively homogenous and crop phenologies are largely synchronous, cropping

122    systems in the tropics are more complex in the sense that planting and harvesting schedules could be

123    substantially different for the same crop (e.g. soybean in Brazil) (Song et al. 2021). Statistics-based

124    phenological metrics derived from time-series of satellite data can capture the salient features of

125    vegetation phenology while maintaining high spatial and temporal data consistency, and thus, provide a

126    unique advantage to large-area vegetation type mapping (DeFries et al. 1995; Hansen et al. 2013; Song et

127    al. 2018). The main objective of this study is to explore the utility of statistical metrics derived from

128    Landsat and MODIS data as well as machine learning algorithms for high-resolution, long-term crop

129    yield mapping in the tropics. Producing long-term spatially explicit yield information is especially

130    imperative in tropical countries, where agricultural production is growing rapidly, causing detrimental

131    impacts to natural environment (Gibbs et al. 2010; Potapov et al. 2022; Song et al. 2018; Zalles et al.

132    2021). We focus on annual soybean yield in Brazil over 2001-2020 in this study.

133    **2.  Data and Methods**

134    **2.1. Study area**

135    Our study area covers the southern hemisphere portion of Brazil. Brazil is the world's leading producer

136    and exporter of soybeans, accounting for more than 35% of global production and about half of the

137    world's total export (FAO 2020). Based on statistics from the Food and Agriculture Organization of the

138    United Nations (FAO), soybean production in Brazil has tripled from 37.9 million tons in 2001 to 114.3

139    million tons in 2019 (FAO 2020). Over the same time period, soybean cultivation area in Brazil increased

140    from 14.0 Mha to 35.9 Mha, and the national average yield increased from 2.71 to 3.18 tons/ha with the

141    maximum yield of 3.39 tons/ha achieved in 2018 (FAO 2020). The dramatic increase in soybean

142    cultivation in Brazil (Figure 1) has directly and indirectly caused widespread natural vegetation loss and

143    cascading environmental impacts in the Amazon, Cerrado and other biomes (Song et al. 2021a; Zalles et

144    al. 2019).

145



146    **Figure 1**. Soybean expansion in Brazil mapped using satellite data. (a) Soybean change during 2001-2010

147    and 2010-2019. For simplicity to visualize, the annual 2001-2019 classification maps are used to create

148    bi-temporal change layers. Landsat mosaic of South America is used as the backdrop in (a), and gray

149    shaded area represents the study area of Brazil. Regional details over two soybean expansion frontiers are

150    shown in (b) Mato Grosso and (c) MaToPiBa (Maranhao, Tocantins, Piaui and Bahia). Reduction in

151    soybean cultivation was observed along the border between Sao Paulo and Minas Gerais, shown in (d).

152

## 2.2. Satellite data and products

154    We used Landsat and MODIS as the main satellite data to derive vegetation characteristics of soybean

155    plants, as they represent the most consistent satellite data records over the past two decades. According to

156    the United States Department of Agriculture (USDA) crop calendars for Brazil, soybeans in Brazil are

157    typically planted in October to December and harvested in March to May

158    (https://ipad.fas.usda.gov/rssiws/al/crop_calendar/br.aspx). In our study, all Landsat and MODIS

159    observations acquired between November 1[st] and April 30[th] of the next year from 2000 to 2019 were

160    processed. The MODIS surface reflectance (SR) data in blue (469 nm), green (555 nm), red (645 nm),

161    near-infrared (NIR, 858 nm), shortwave infrared (SWIR, 1640 nm and 2130 nm) and thermal (11,030 nm)

162    wavelengths were obtained as 16-day composites from the MOD44C product, same as the MOD09GA,

163    MOD09GQ and MODTBGA v006 products (Vermote and Wolfe 2015). Landsat images acquired by the

164    Thematic Mapper (TM), Enhanced Thematic Mapper Plus (ETM+), and Operational Land Imager (OLI),

165    with blue, green, red, NIR, and SWIR bands, were converted from top-of-atmosphere reflectance to

166    normalized surface reflectance (NSR) through an automated data processing system (Potapov et al. 2020).

167    Using MODIS SR as normalization target, the system corrected atmospheric and anisotropic effects of

168    Landsat after at-sensor radiance calculation, cloud, shadow and haze masking. The Landsat NSR, from all

169    sensors, was then processed to 16-day composites consistent with the MODIS product. Both Landsat

170    NSR and MODIS SR 16-day time-series were used to create seasonal phenological metrics, including

171    NDVI, EVI, normalized difference water index (NDWI) and other band ratio indices (Table 1). A

172    complete description of Landsat data processing and the freely available software tools to generate

173    phenological metrics is provided in Potapov et al. (2020).

174    **Table 1**. Input features for modeling and mapping soybean yield in Brazil. Please see Supplementary

175    Information for the complete list of variables.

| Category | Input Features | N |
|---|---|---|
| Landsat-based | Seasonal vegetation phenological metrics derived from Blue, Green, Red, NIR, SWIR1, SWIR2 and thermal bands | 50 |
| MODIS-based | Seasonal vegetation phenological metrics derived from Blue, Green, Red, NIR, SWIR1, SWIR2 and thermal bands | 24 |
| Topographic | DEM and Slope | 2 |
| Climate | Long-term (1971-2000 average) climate data, monthly (October to May) TMP (mean 2 m temperature), DTR (diurnal 2 m temperature range), PRE (precipitation rate), VAP (vapor pressure), WET (wet days), CLD (cloud cover), TMN (minimum 2 m temperature), TMX (maximum 2 m temperature) and PET (potential evapotranspiration) | 72 |
| Weather | Annual (2000 through 2019) in-season weather data, monthly (October to May) TMP, DTR, PRE, VAP, WET, CLD, TMN, TMX and PET | 72 |
| Soil | Water storage capacity, topsoil and subsoil bulk density, cation exchange capacity of the clay fraction in the topsoil and subsoil, topsoil and subsoil clay, sand and silt fractions, topsoil and subsoil pH, and area weighted topsoil and subsoil carbon content | 15 |

176

177    We used a recently developed 30 m resolution (0.00025° × 0.00025°), annual, 2001-2019 soybean

178    classification map product (Song et al. 2021a) as masks to constrain the yield modeling and mapping to

179    identified soybean pixels (Figure 1). For simplicity and consistent with the soybean classification map

180    product, in this study we refer to a cropping year by the harvest year. For example, year 2001 indicates

181    the 2000/01 cropping year. The soybean classification product was developed using the above Landsat

182    and MODIS data as input in addition to 30 m resolution topographic features from the Shuttle Radar

183    Topography Mission (SRTM) data. Continentally distributed field observations collected over three years

184    (2017, 2018 and 2019) were used as training to calibrate a multi-year bagged decision tree model for

185     soybean classification. The overall accuracy of the soybean classification maps for the years of 2017,

186     2018, and 2019, where we had probability field sample for validation, was 96%, 94% and 96%,

187     respectively, with high and balanced producer's and user's accuracies (Song et al. 2021a).

188     **2.3. Climate and weather data**

189     Monthly climate and weather covariates were obtained from the Climatic Research Unit gridded Time

190     Series (CRU TS) version 4.04 dataset (Harris et al. 2020). The variables included TMP (mean 2 m

191     temperature), DTR (diurnal 2 m temperature range), PRE (precipitation rate), VAP (vapor pressure),

192     WET (wet days), CLD (cloud cover), TMN (minimum 2 m temperature), TMX (maximum 2 m

193     temperature) and PET (potential evapotranspiration) at a spatial resolution of $0.5° \times 0.5°$. We calculated

194     monthly average values from 1971 to 2000 for the months from October to May to represent long-term

195     climatology. For each year between 2000 to 2019, we directly used the monthly values for the months

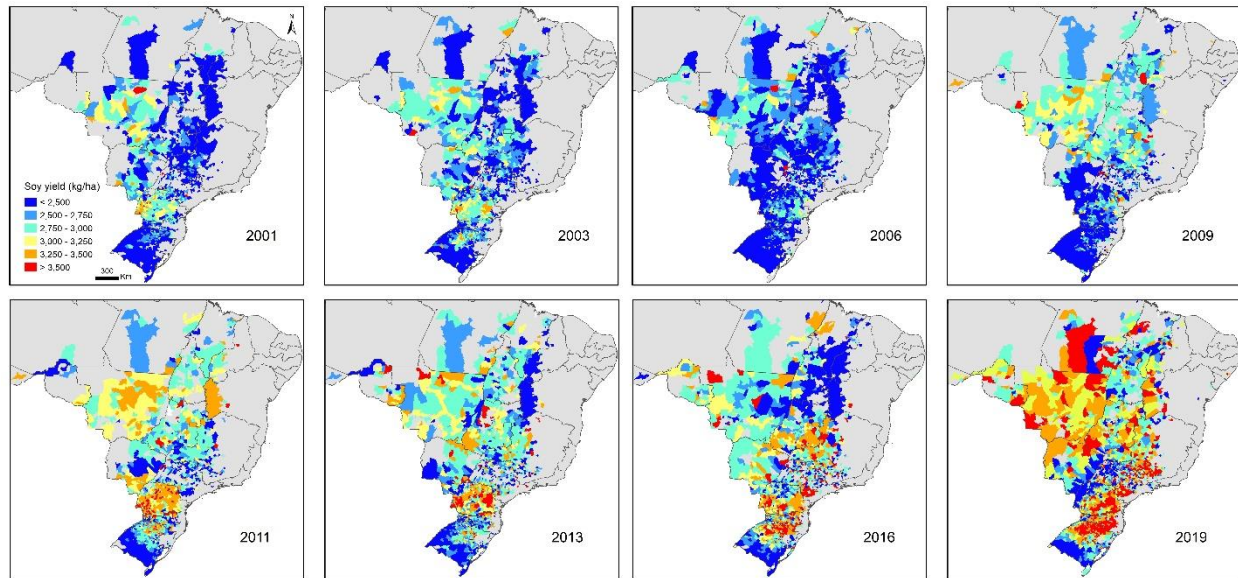196     from October to May to represent in-season weather (Table 1).

197     **2.4. Soil data**

198     The Regridded Harmonized World Soil Database v1.2 at $0.05° \times 0.05°$ spatial resolution

199     (FAO/IIASA/ISRIC/ISSCAS/JRC 2012; Wieder et al. 2014) were obtained and processed similar to the

200     climate and weather data. The soil variables included available water storage capacity, topsoil (0-30 cm)

201     and subsoil (30-100 cm) bulk density, cation exchange capacity of the clay fraction in the topsoil and

202     subsoil, topsoil and subsoil clay, sand and silt fractions, topsoil and subsoil pH, and area weighted topsoil

203     and subsoil carbon content (Table 1).

204     **2.5. Municipal yield statistics**

205     We obtained soybean yield statistics at the municipality scale for every year between 2001 and 2019 from

206     the Brazilian Institute of Geography and Statistics (IBGE) Municipal Agricultural Production database

207     (https://sidra.ibge.gov.br/). The size of the municipalities where soybeans are cultivated varies widely

208     from south (small) to north (large), with a median size of approximately 48 Kha, the first quantile of 22
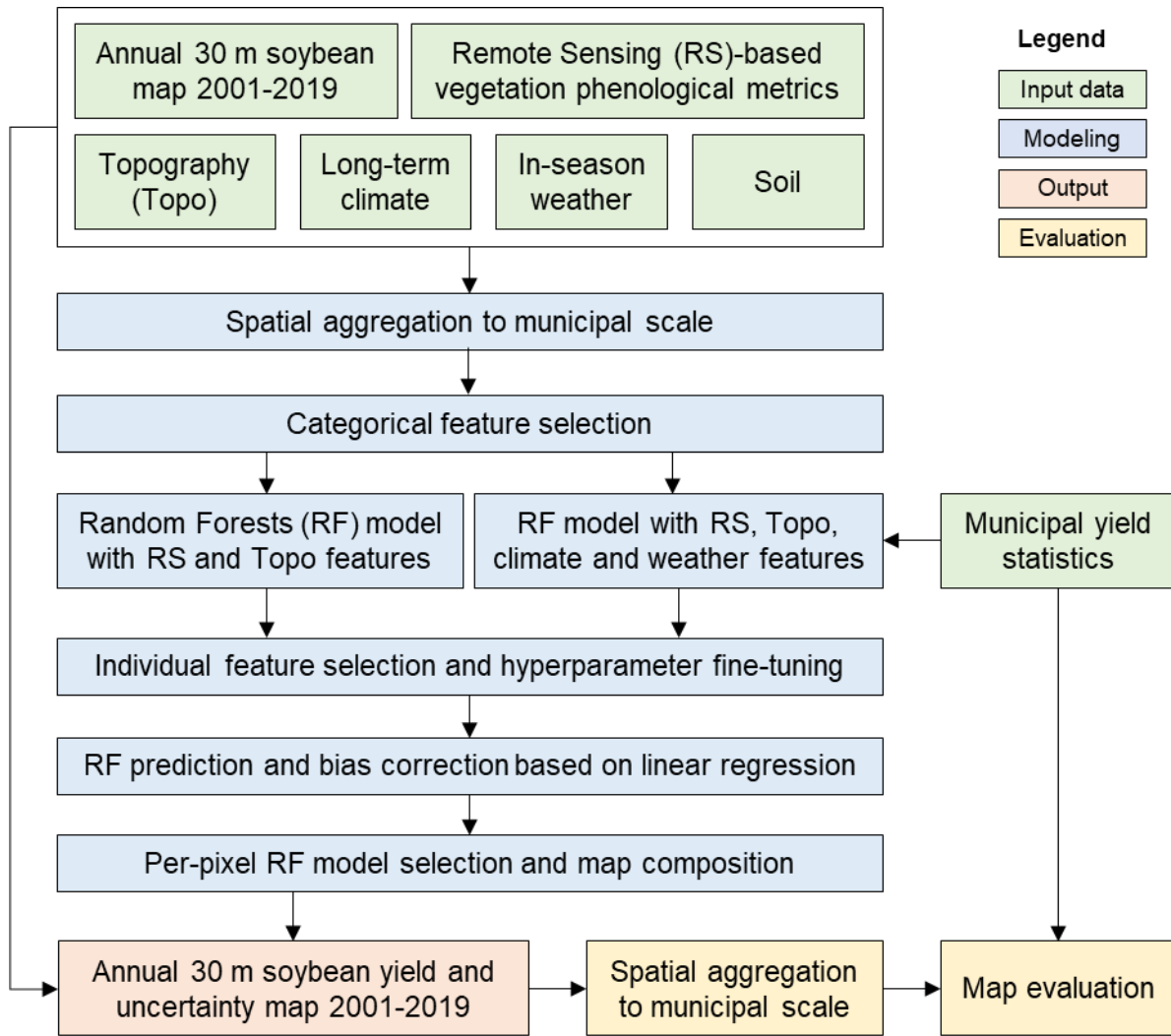
209    Kha and the third quantile of 135 Kha. These yield statistics were used as reference data for training and

210    evaluation (Figure 2).



211

**Figure 2**. Municipality-level yield statistics from the Brazilian Institute of Geography and Statistics
213    (IBGE) were used as reference for modeling and mapping soy yield.

214

## 2.6. Modeling yield

216    The overall workflow of modeling and mapping soybean yield is presented in Figure 3. Major steps

217    include spatial aggregation of remote sensing (RS)-based vegetation phenological metrics, topographic

218    (topo) features, climate, weather, and soil variables to municipal scale, categorical feature selection,

219    random forests (RF) (Breiman 2001) model training, RF prediction, bias correction, per-pixel RF model

220    selection and composition, and map evaluation. Details of each step are described as follows.

221

**Figure 3**. Overall workflow of mapping annual soybean yield 2001-2019 using satellite data, climate,

weather, soil and topography data, municipality statistics, random forests and linear regression models.

Two random forests models were trained and implemented with more details reported in the text.

The $0.5° \times 0.5°$ climate and weather data, and the $0.05° \times 0.05°$ soil data were first resampled using

nearest resampling to $0.00025° \times 0.00025°$ to match the spatial resolution of the soybean classification

map, remote sensing data and topographic features. With the annual soybean classification map as a

mask, we aggregated these input datasets to municipal scale by taking the average value over soybean

pixels in each municipality. The spatial aggregation step was conducted for every year independently

231    between 2001 and 2019. To remove the non-soybean and low-soybean municipalities, we selected the

232    municipalities with annual soybean pixels ≥ 50,000, resulting in a total of 15,784 municipalities across the

233    19-year period. These municipalities contained 95% of all mapped soybean pixels over the study period.

234    To investigate the relative utilities of these multi-source, multi-resolution input datasets for yield

235    modeling, we conducted three progressive experiments using categorical feature selection. Specifically,

236    we built three random forests models with (1) RS and topo features as input, (2) RS, topo, climate and

237    weather features as input, and (3) RS, topo, climate, weather and soil features as input. Performance of

238    model #1 represents the utility of RS and topo features to model yield. Improved performance of model

239    #2 over model #1 would represent the value of weather and climate data. Likewise, improved

240    performance of model #3 over model #2 would represent the value of the soil variables.

241    Municipal yield statistics were used as reference for all three models. For each model, we randomly

242    selected 80% municipalities as training (n = 12,649) and the remaining 20% was reserved for independent

243    test (n = 3,135), with both training and test data covering all 19 years. We calculated root-mean-square-

244    error (RMSE), mean bias error (MBE), mean absolute error (MAE), and $r^2$ using both training and test

245    data for all three models. To further enhance the robustness of the model evaluation and to eliminate

246    potential bias from a particular realization of sampling, we implemented a Monte Carlo method and

247    repeated the random training/test split, model training and evaluation 100 times. The final model

248    performance was represented using box plots of RMSE, MBE, MAE and $r^2$ of the 100 runs.

249    In addition to model evaluation with 20% withheld test data, we also conducted the leave-one-year-out

250    model assessment. For every year between 2001 and 2019, we used 18-years of data to calibrate the

251    random forests models and used the model to predict over the left-out year. For the left-out year, we

252    compared the predicted yield with reference statistics and calculated error metrics.

253    Our model assessment revealed that climate and weather variables significantly improved model

254    performance, but soil variables did not further improve model performance (more details are provided in

255    the Results and Discussion sections). Therefore, the model with RS, topo, climate and weather variables

256    as input (i.e. model #2) was selected as the primary model for yield estimation. However, due to the

257    coarse spatial resolution ($0.5° \times 0.5°$) of the climate and weather data, spatial grid patterns were noticed in

258    some regions. To remove these artifacts, we implemented model #1 (RS and topo features as input) as a

259    secondary model, and results of the two models were combined (see more details below).

260    To improve computational efficiency, we conducted individual feature selection for both models. For

261    each RF model, we trained the model using all features as input, ranked each feature and selected the top

262    features with a cumulative importance of greater than 95%. We also constructed a correlation matrix of

263    the features and removed those less important features that had a correlation coefficient of greater than

264    0.95 with the more important ones. Error metrics were calculated for all as well as selected features to

265    demonstrate the comparable performance of trained models. We implemented the random forest classifier

266    function in the sklearn package in python. The RF parameters fine-tuned included n_estimators (number

267    of trees), max_features (number of features to consider at every split), max_depth (maximum number of

268    levels in a tree), min_samples_split (minimum number of samples required to split a node),

269    min_samples_leaf (minimum number of samples required at each leaf node). We applied a randomized

270    search on hyper-parameters followed by a grid search to determine the exact values for these parameters.

271    The immediate output of the two RF models include predicted soybean yield, represented as the mean

272    value of all trees in the forest, and associated uncertainty, represented as the standard deviation of all trees

273    in the forest. For continuous variables, random forests could generate underestimation at the high-end of

274    the variable and overestimation at the low-end of the variable because of the effect of "regression to the

275    mean" (Huang et al. 2016; Zhang and Lu 2012). Such is the case for our yield modeling in this study. To

276    correct these systematic biases, we followed Zhang and Lu (2012) and Huang et al. (2016), and applied

277    linear regression using the municipal yield statistics as the dependent variable and the RF-predicted yield

278    as the independent variable. The derived linear equation was subsequently applied to the adjust the RF-

279    predicted yield and uncertainty.

280    We implemented the two calibrated random forest models (models #1 and #2) and their associated linear

281    regressions independently using the annual input datasets. The outputs were two sets of 30 m resolution

282    soybean yield and uncertainty maps for every year between 2001 and 2019. We created a final soybean

283    yield and uncertainty map for every year through per-pixel composition, where, for every pixel, the

284    soybean yield and associated uncertainty were selected from the model with a smaller uncertainty.

285    **2.7. Yield map evaluation**

286    We evaluated the quality of the annual, 30 m resolution soybean yield maps at the municipal scale.

287    Average yield was derived from the maps, and compared to municipal yield statistics as reference. We

288    computed the difference of the two datasets and constructed a histogram. We calculated RMSE, MAE,

289    MBE, and $r^2$, and created scatter plots using the 19 years of data. We also calculated these error metrics

290    for every year to evaluate the temporal consistency of the yield map time series.
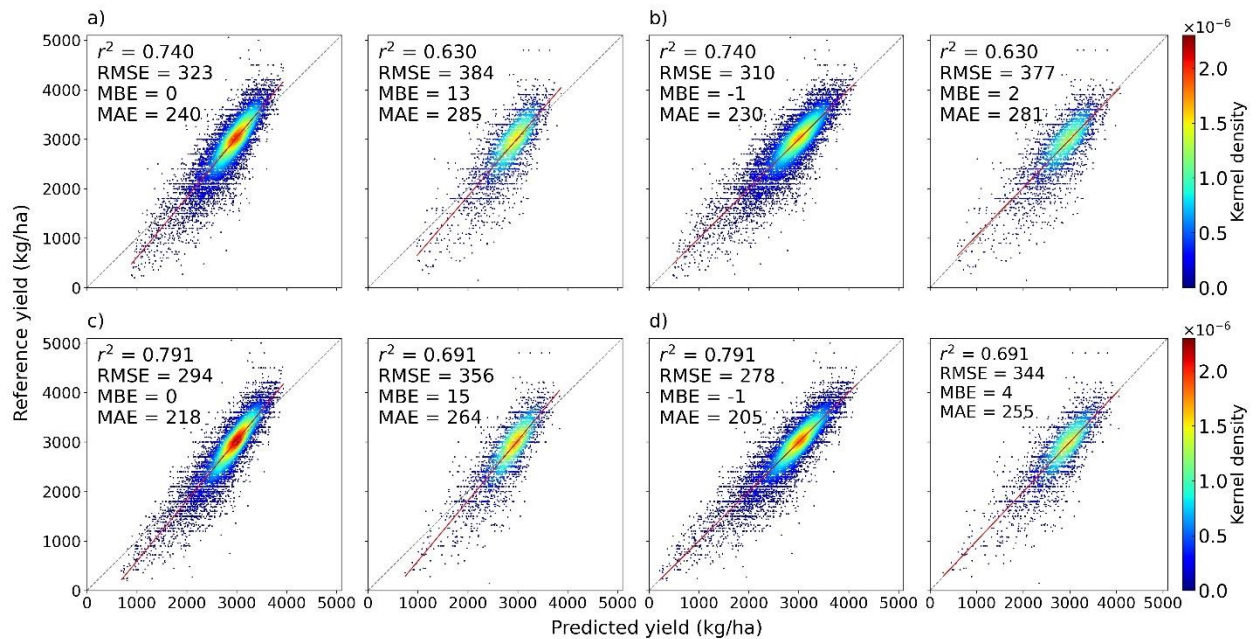
291    **3.   Results**

292    **3.1. Model selection and performance**

293    Using remote sensing-based vegetation phenological metrics and topographic features as input to random

294    forests (model #1) produced an $r^2$ of 0.74, an RMSE of 323 kg/ha, an MBE of 0 kg/ha and a MAE of 240

295    kg/ha for training data. Compared to the 2001-2019 national average yield of 2,869 kg/ha, this RMSE

296    represents 11% error. Adding climate and weather variables to input (model #2) significantly improved

297    model performance, as represented by the increase in $r^2$ and reduction in RMSE and MAE, for both

298    training and test data. The improved model had an $r^2$ of 0.79, an RMSE of 294 kg/ha, an MBE of 0 kg/ha

299    and a MAE of 218 kg/ha for training data, and an $r^2$ of 0.69, an RMSE of 356 kg/ha, an MBE of 15 kg/ha

300    and a MAE of 264 kg/ha for test data. Adding soil variables to input (model #3) showed little to no value

301    in further improving model performance. Therefore, we discarded model #3 and implemented model #1

302    and #2 in this study. Both model #1 and #2 were chosen because although climate and weather data

303    demonstrated considerable utility in modeling soybean yield, their coarse spatial resolution ($0.5° \times 0.5°$)

304  caused apparent grid patterns when the model was applied to 30 meter spatial resolution, whereas model

305  #1 generated spatially coherent results. Moreover, individual feature selection not only improved

306  computational efficiency but also improved model accuracy. Consistent for all model categories, there

307  remained some differences between training and test, indicating potential overfitting of the models. This

308  was likely due to the lack of high-quality soil data and other important agricultural management variables

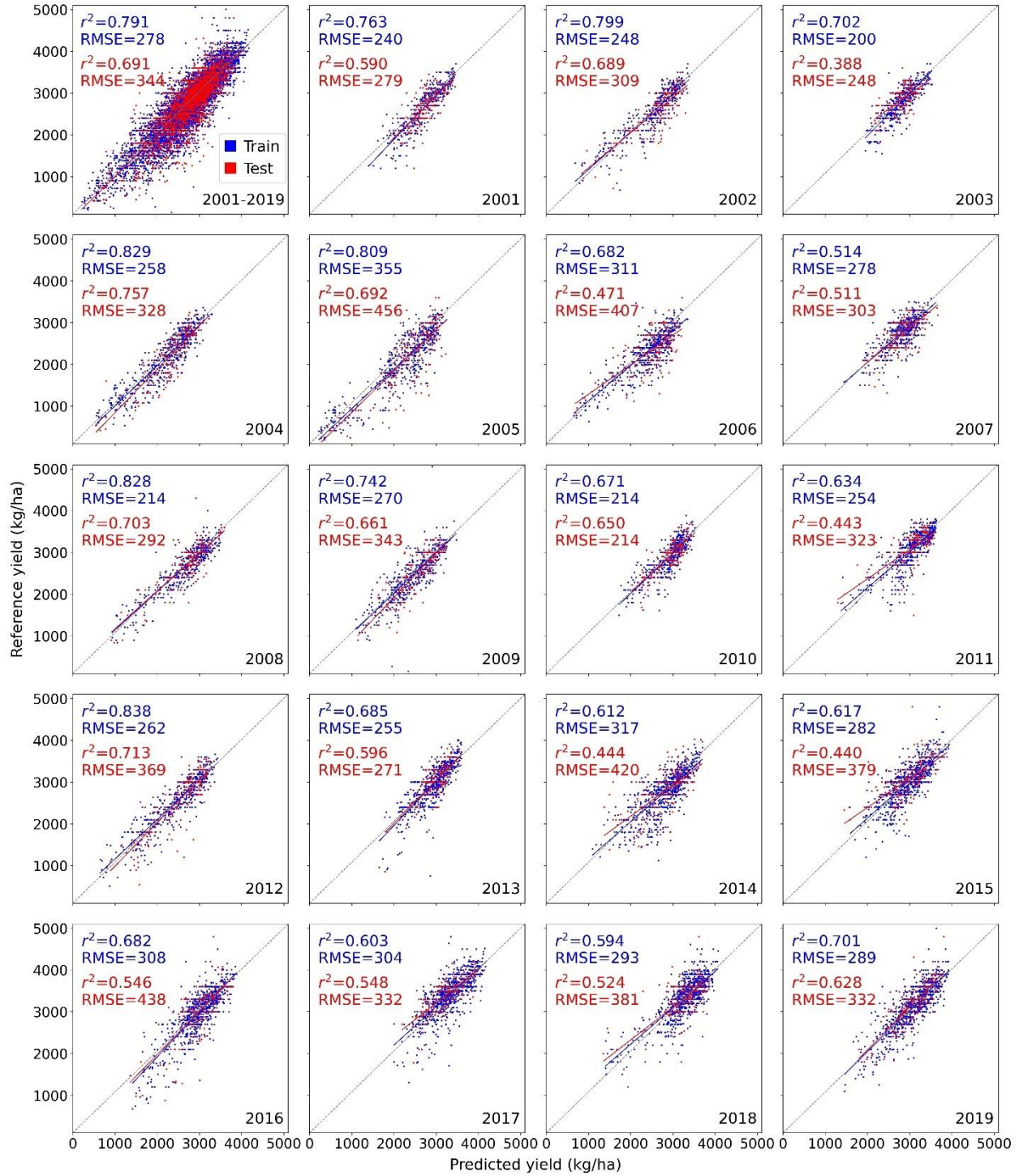309  (e.g. fertilizer use) in the model (please see more details in the Discussion section).

310  Predicted yield from random forests models were highly consistent with reference yield from municipal

311  statistics (Figure 4). However, the direct outputs of the random forests models under-estimated yield at

312  the high end and over-estimated yield at the low end (Figure 4a and 4c). Applying a linear regression

313  successfully corrected these systematic biases for both models (Figure 4b and 4d). Moreover, the overall

314  model performance was also slightly improved, as demonstrated by the reduction in RMSE and MAE for

315  both training and test results. For instance, the training accuracy in terms of RMSE was reduced from 294

316  to 278 kg/ha and the test accuracy was improved from 356 to 344 kg/ha for model #2 after bias

317  adjustment (Figure 4a vs 4b).



318

319 **Figure 4**. Performance of yield models before and after systematic bias adjustment using linear

320 regression. a) Random forests (RF)-predicted soybean yield against reference yield from municipal

321 statistics. Input data for RF include remote sensing, topographic features, climate and weather variables.

322 The left panel is density scatter plots using training data and the right panel is density scatter plots of

323 independent test data. The red lines on both panels represent the linear regression line. b) Same as a), but

324 a linear regression was applied to adjust bias in RF outputs. c) RF-predicted soybean yield against

325 reference yield. Input data for RF only include remote sensing and topographic features. d) Same as c),

326 but after linear bias adjustment.

327

328 Although the model was trained using all 19-years of data as input, evaluation of model performance at

329 the annual time scale revealed consistent model performance across all 19 years (Figure 5). Based on the

330 withheld test data, the 19-year overall RMSE was 344 kg/ha and the $r^2$ was 0.69. The RMSE represents

331 12% error relative to long-term yield mean. The annual RMSE values ranged from 214 kg/ha in 2010 to

332 456 kg/ha in 2005, and the annual $r^2$ values ranged from 0.39 in 2003 to 0.76 in 2004. No significant

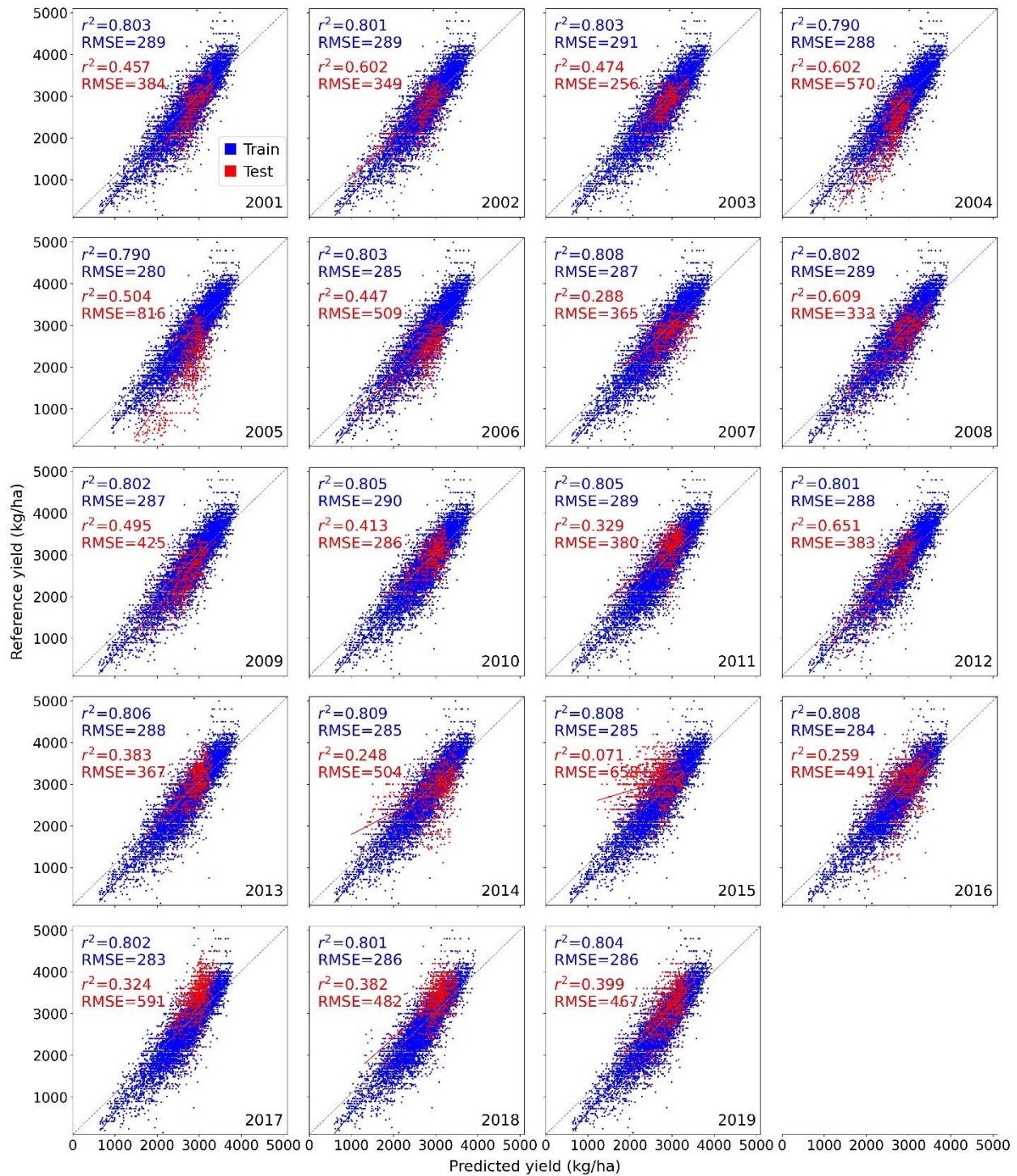333 systematic bias was observed for any of the years (Figure 5).

334 The leave-one-year-out model assessment revealed that the yield models performed well for most of the

335 19 years, but performed relatively poorly for 2005 and 2015 with notably higher RMSE and lower $r^2$,

336 respectively (Figure 6). The RMSE of the leave-one-year-out assessment ranged from 259 kg/ha to 816

337 kg/ha. These results are in general comparable to regional studies of satellite-based soybean yield

338 mapping in the Midwest of the United States (Lobell et al. 2015) and Southern Brazil (Schwalbert et al.

339 2020). Both 2005 and 2015 did not show notable performance deficiency when data of the two years were

340 included in training (Figure 5). Comparison between annual accuracies of the two model assessments

341 (Figures 5 and 6) suggests that model trained with long time series of data generally perform well for

342 unseen years. The comparison also highlights the significance of including both good and poor harvesting

343 years in training for enhancing the temporal generalization and predictive capability of trained models.

344

**Figure 5**. Performance of yield model at an annual time scale. X-axis represents model-predicted yield, and y-axis represents reference yield from municipal statistics. The top-left scatter plot is a combination of the two scatter plots in Figure 5b. Scatter plots are made using training data and withheld test data. Input data for model include remote sensing, topographic features, climate and weather variables.
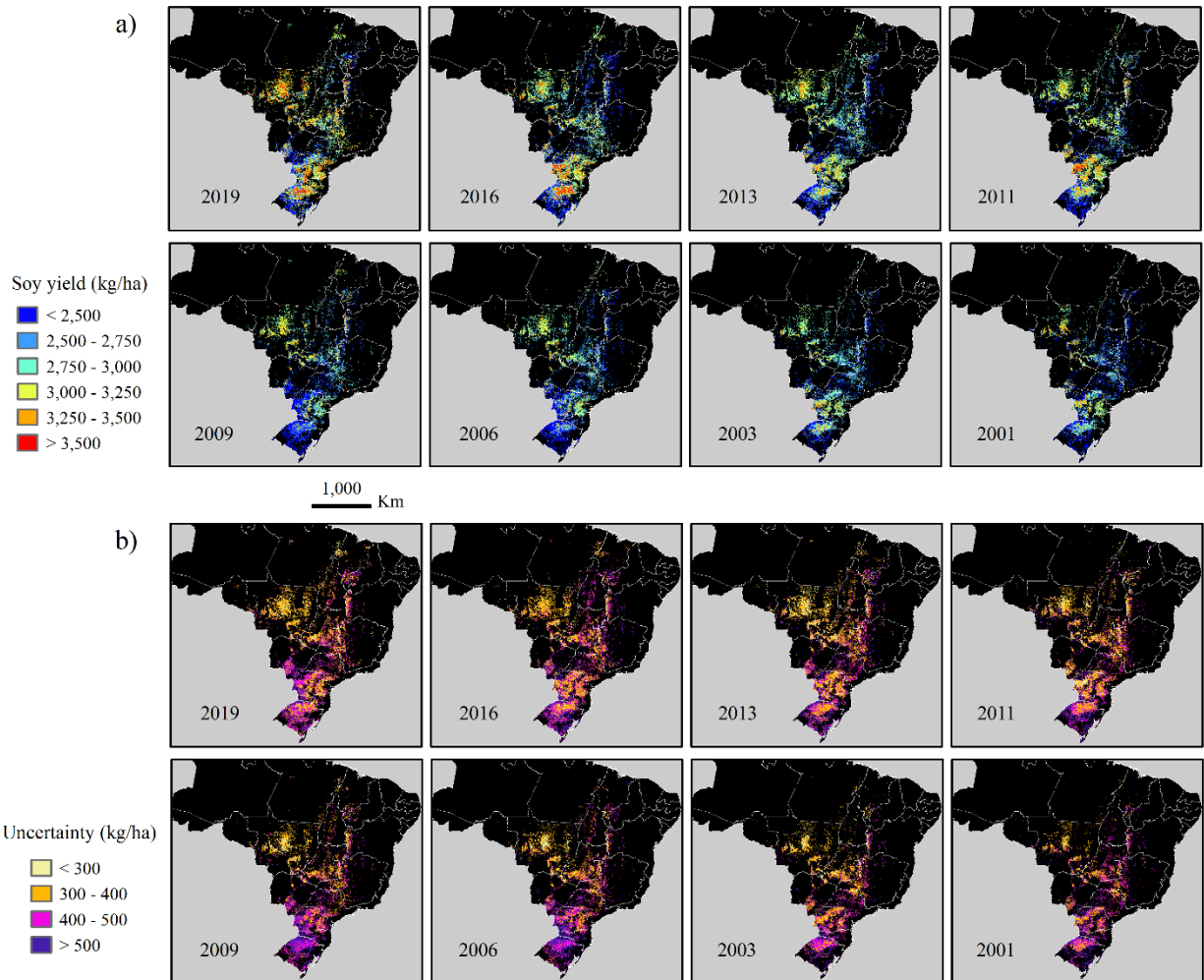
350



351

**Figure 6**. Leave-one-year-out model assessment. For each year between 2001 and 2019, 18-years of data were used to training the model (blue dots and text), which was used to predict over the left-out year.

354    Municipal statistics of the left-out year were used as reference to evaluate the model performance (red

355    dots and text).

356

357    **3.2. Annual soybean yield and uncertainty maps**
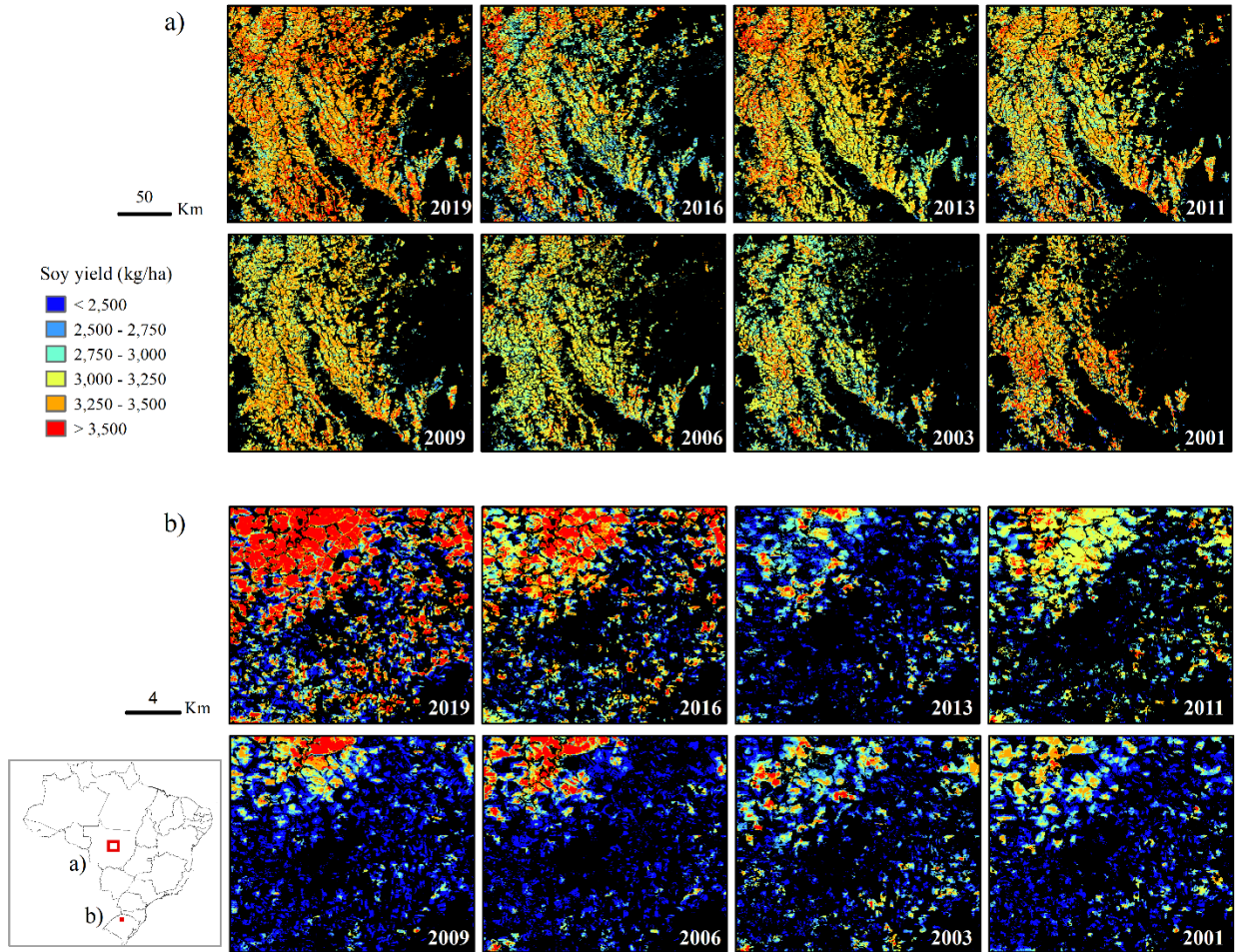
358    Implementing the calibrated random forests and linear regression models at 30 m spatial resolution

359    generated spatially and temporally coherent soybean yield distributions across Brazil from 2001 to 2019

360    (Figure 7a). Considerable spatial heterogeneity in soybean yields was observed across the country. In

361    2001, the highest soybean yield regions included central Mato Grosso and western Parana (also see Figure

362    2a), and the lowest yield regions included Rio Grande do Sul, eastern Goias, western Minas Gerais, and

363    western Bahia. Increase in soybean yield was found in many regions, most notably in northern Rio

364    Grande do Sul and western Bahia (also see Figure 2b). Soybeans in Mato Grosso experienced not only a

365    substantial area expansion but also considerable yield growth. Per-pixel uncertainty of soybean yields

366    (Figure 7b) showed that the uncertainty estimates were mostly between 300 kg/ha to 500 kg/ha.

367    Moreover, the uncertainty distribution varied both spatially and temporally, with the south region (e.g.

368    Rio Grande do Sul) appeared to have slightly higher uncertainties than center west (e.g. Mato Grosso).

369

370 **Figure 7.** Annual soybean yield and uncertainty maps for selected years over Brazil. Yield and

371 uncertainty maps were produced at 30 m spatial resolution and averaged to 1 km for the purpose of

372 display. Regional details at 30 m resolution are shown in Figure 8.

373

374 The annual, 30 m resolution maps revealed field-level heterogeneity in soybean yields (Figure 8). Large

375 contiguous soybean fields in central Mato Grosso have moderate-to-high yield and small variations

376 between fields (Figure 8a), whereas smaller fragmented fields in Rio Grande do Sul show much larger

377 variations (Figure 8b). Over the past 19 years, soybean yields in central Mato Grosso experienced an

378 overall increase in most fields, whereas in Rio Grande do Sul, larger fields appeared to have relatively

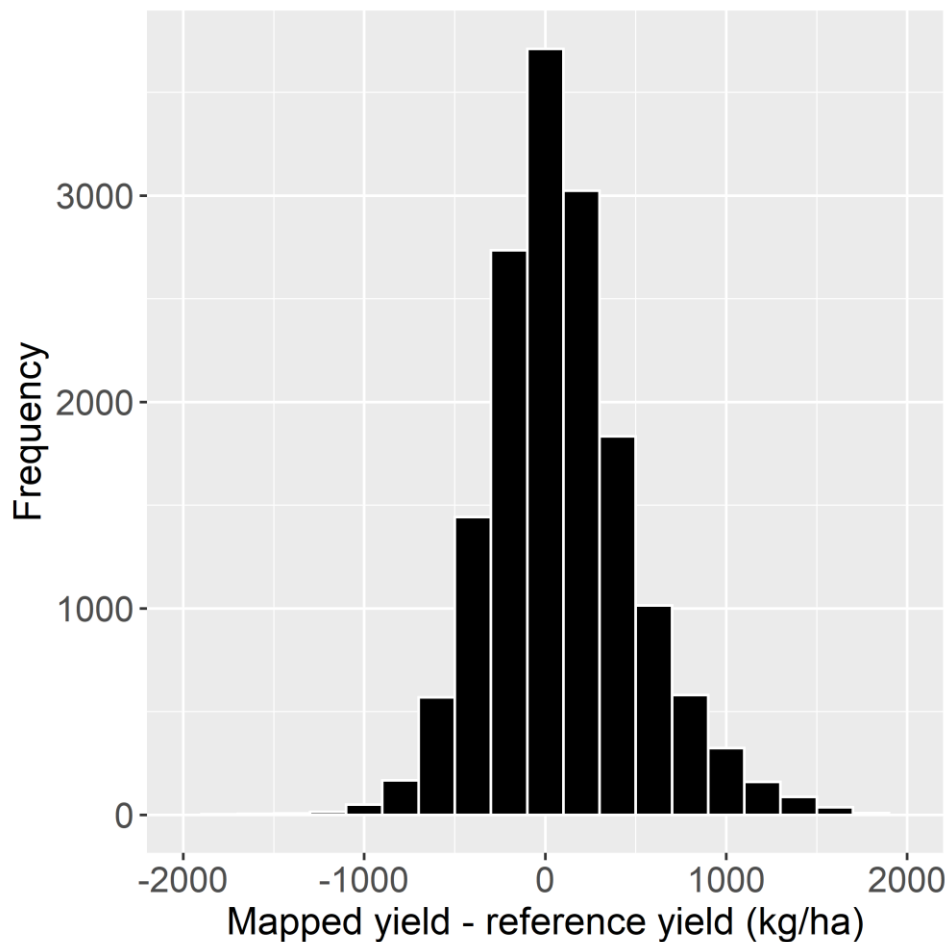379 greater yield growth than smaller fields (Figure 8b).

**Figure 8.** Spatial and temporal details of soybean yield at 30 m resolution in two selected regions: **a)** central Mato Grosso and **b**) northern Rio Grande do Sul. Field-level yield heterogeneity is revealed by the time series of high-resolution maps.
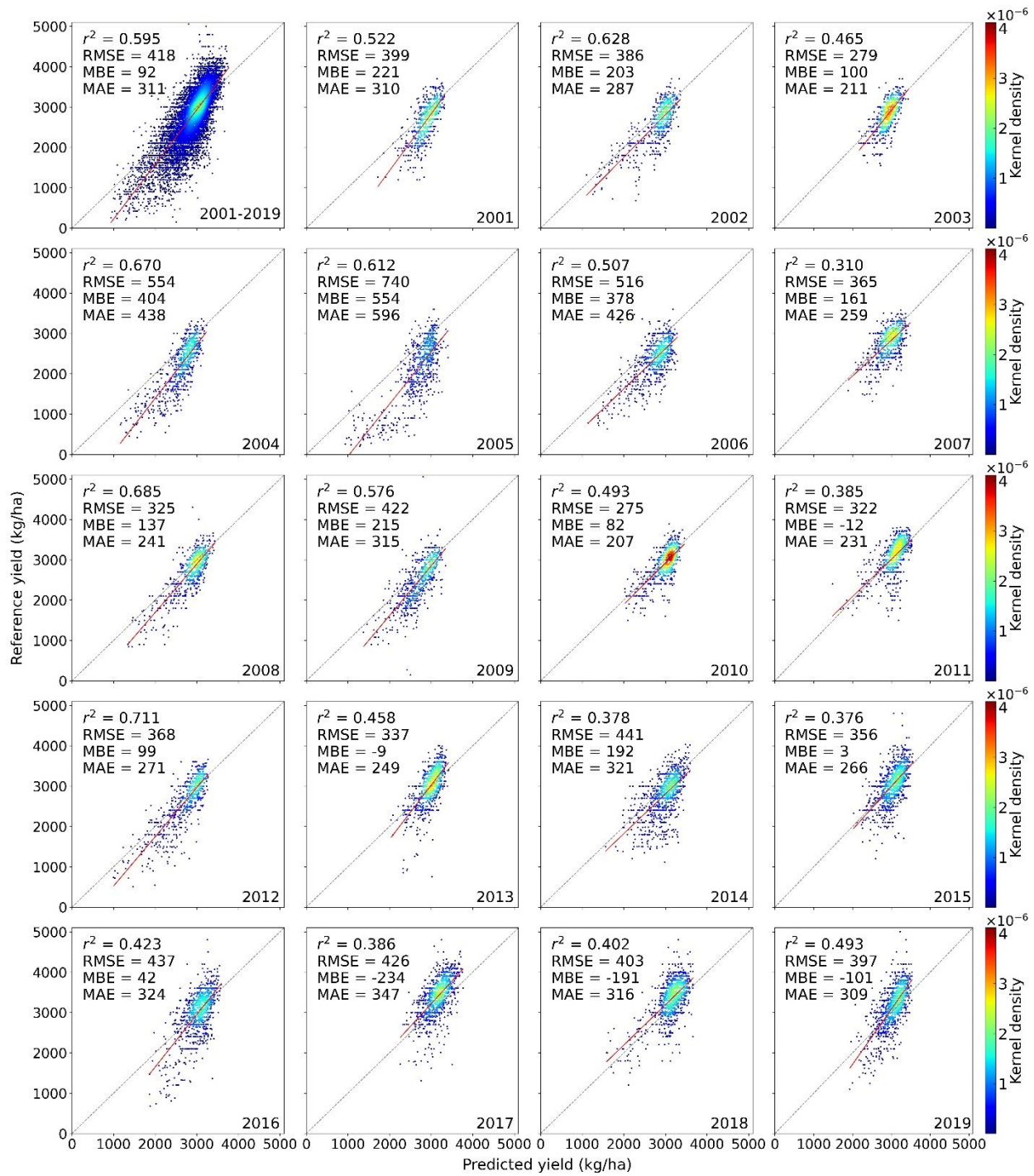
### 3.3. Map evaluation

The annual 30 m soybean yield maps were aggregated to municipal scale for a quantitative quality assessment. Compared to the reference data from official statistics, the yield map product had an overall RMSE of 418 kg/ha, a MAE of 311 kg/ha, an MBE of 92 kg/ha, and an $r^2$ of 0.60. Compared to the 2001-2019 national average yield of 2,869 kg/ha, the RMSE represents 15% error. These error metrics were all slightly worse than the model performance, with the RMSE about 20% higher (compared to 344 kg/ha; see detailed numbers of other error metrics in Figure 4). An overall slight positive bias was noted (mean bias of 92 kg/ha or 3% error compared to long-term average yield, Figure 9). Moreover, systematic

392    underestimation was still noticed at the high end of yield and overestimation at the low end of yield

393    (Figure 10), although a linear regression successfully corrected model bias at the training stage at the

394    municipal level (Figure 4). At the annual time scale, the map accuracy was comparable to model

395    performance for the majority of the 19 years (Figure 10). The comparison between model performance

396    and map quality assessment suggested that uncertainties at the 30 m pixel scale were larger than those at

397    the aggregated municipal scale, highlighting a general multi-scale issue in the applications of regression-

398    based machine learning algorithms in remote sensing.



399

400    **Figure 9.** Histogram of the difference between predicted yield and reference yield at the municipal level

401    between 2001 and 2019 (n=15,784) indicating a slight positive bias in the predicted yield.
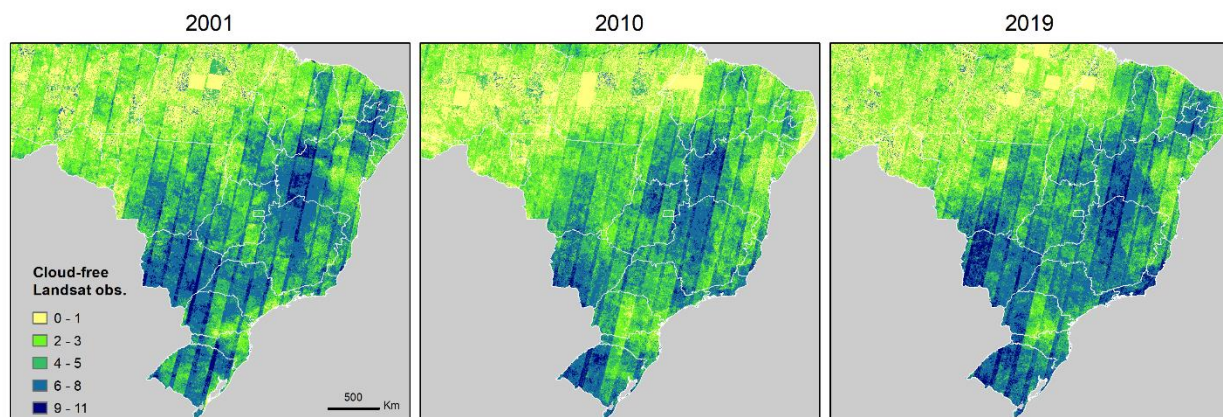
**Figure 10.** Quality assessment of 30 m soybean yield maps for every year between 2001 and 2019. The annual maps were averaged to the municipal scale to derive predicted yields (x-axis). Reference yields (y-axis) are official statistics.

407 **4. Discussion**

408 **4.1. Uncertainty sources for yield modeling**
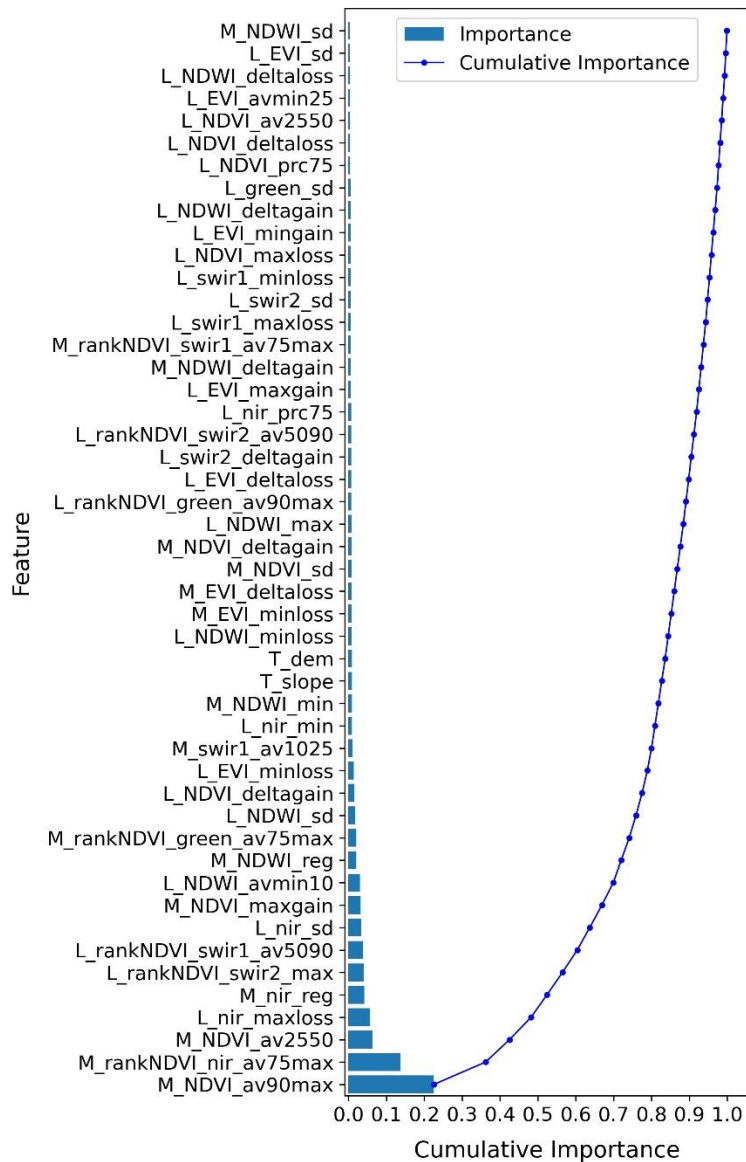
409 Model performance and the quality of the annual yield maps are influced by a number of factors,

410 including the temporal density of satellite observations, the coarse spatial resolution and uncertainties of

411 climate and weather variables, lack of up-to-date soil measurements, unknown uncertainties in the official

412 statsitics, lack of field-level reference data, missclassifications in the annual soybean masks, and the muti-

413 scale modeling and prediction procedure. The impacts of these factors are discussed in detail as follows.

414 Depending on the type of cultivar, environmental conditions and agricultural management practices,

415 soybean plants take 90 to 150 days from planting to maturity. During this short growing window,

416 vegetation cover in the field experiences rapid transitions from bare ground to nearly closed canopy and

417 to bare ground again. Such phenological dynamics require dense time-series data to capture the key

418 growth stages that are critical to crop biomass accumulation and yield formation. Studies have

419 demonstrated that the peak growing period in vegetation index is most important for modeling yield for

420 wheat, corn and soybeans (Becker-Reshef et al. 2010; Johnson 2014). In addition, natural disasters during

421 or after the seed-filling stage can cause severe yield reduction (Hosseini et al. 2020). In this study, we

422 used MODIS and Landsat as the main remote sensing data source. Due to the sparse temporal interval of

423 Landsat, cloud-free Landsat observations vary considerably in space and time (Figure 11).

424

425  **Figure 11.** Cloud-free Landsat observations between November 1st and April 30th in selected years over

426  Brazil.

427

428  On the other hand, daily MODIS acquisitions are more robust to cloud contamination. Indeed, the

429  important features identified by random forests include many MODIS-based spectral features. The most

430  important feature of the random forests model (model #1) was "M_NDVI_av90max", which represented

431  the average value of the 90th percentile and maximum NDVI (i.e. peak NDVI) derived from MODIS

432  (Figure 12). The second and third most important features were MODIS-based peak-season NIR

433  reflectance and middle-season NDVI, respectively. These top three features accounted for >40% of

434  cumulative feature importance (Figure 12). Another inherent factor that enabled MODIS to be an efficient

435  sensor for modeling soybean yield is the large field size in Brazil (Fritz et al. 2015). The feature ranking

436  analysis suggested that improving the temporal density of high spatial resolution satellite data, such as the

437  Harmonized Landsat and Sentinel-2 product (Claverie et al. 2018), may improve yield mapping at the

438  field scale. Further research is also needed to investigate the utility of other freely available satellite data,

439  particularly radar data (e.g. Sentinel 1) for yield estimation, as radar data can provide complementary

440  infromation to optical data for crop monitoring (Song et al. 2021b; Veloso et al. 2017) in addition to their

441  all-weather data acquisition.

442

**Figure 12.** Cumulative feature importance for the random forests-based soybean yield modeling using MODIS and Landsat phenological metrics as input. Features with a prefix of "M*" represents MODIS-based metrics, features with a prefix of "L*" represents Landsat-based metrics, and features with a prefix of "T*" represents topographic variables. "av" stands for "average". The metrics are sorted from high to low along the vertical axis from bottom to top. Please see supplementary Table S1 for more explanation of metric names.

449

450   Our study explicitly demonstrated the value of climate and weather data for modeling crop yield. For the

451   trained random forests model with all the features as input, climate and weather variables accounted for

452   36% of the total feature importance (Table 2). Compared to the models with only remote sensing data as

453   input, adding climate and weather variables reduced RMSE by about 7 to 9%, and the improvement was

454   statistically significant. However, adding coarse-resolution climate and weather variables could also

455   introduce undesirable artifacts. By constructing two models and through per-pixel composition of model

456   outputs, our strategy effectively combined the advantages of the two respective models. For any given

457   year, the primary model (i.e. the one with climate and weather variables as input) was chosen for the

458   majority of soybean growing regions of the country, while the secondary model (i.e. the one without

459   climate and weather variables) was selected only for some clustered regions (Figure 13). This data-driven

460   approach relied on the explicit uncertainty outputs associated with predictions of random forests, and the

461   composited map had minimum uncertainties from the multi-model ensemble. Future research will

462   evaluate the uncertainty of climate and weather variables to yield estimation, and incorporate higher-

463   resolution weather dataset for improved yield estimation, e.g. the Climate Hazards Group Infrared

464   Precipitation with Stations (CHIRPS) precipitation data (Funk et al. 2015).

465

466   **Table 2**. Importance of the five categories of input variables in random forests model for soybean yield

467   prediction. Details of the variables are listed in Table 1. The total importance of all variables within each

468   category was calculated and reported.

| Category of variables | Importance in random forests model |
| --- | --- |
| Landsat-based | 0.1883 |
| MODIS-based | 0.4371 |
| Climate | 0.1037 |
| Weather | 0.2539 |
| Topographic | 0.0041 |
| Soil | 0.0128 |

469

470

**Figure 13.** Maps of random forests models chosen for predicting annual soybean yield. The model with
climate and weather varaibles as input was more accurate and was used in the majority of the soybean
growing regions of the country in every year.

The lack of contribution by soil variables to soybean yield modeling was likely because the soil data were
outdated. Soil characteristics and topography are strong determinant of cropland suitability (Ishikawa and
Yamazaki 2021). We used the Harmonized World Soil Database (HWSD) in this study, which was
complied from multiple data sources (FAO/IIASA/ISRIC/ISSCAS/JRC 2012). The data source for Brazil
was the Soil and Terrain database for Latin America and the Caribbean, at the scale of 1:5 million and
released in 1998. Therefore, HWSD represents the soil conditions in Brazil before 1998. From 2000 to
2019, soybean cultivation area in Brazil nearly tripled, and new soybean fields were mostly converted
from pasture and forests (Song et al. 2021a). The conversion process involves removal of surface
vegetation and extraction of the root systems. Subsequently, soil prerparation is critical for cultivating
soybeans on the newly converted land. In the Cerrado, the largest soybean growing biome in Brazil, the
native soil condition is poor for crop production. Most of the soils in the Cerrado are highly weathered
Oxisols and Ultisols, with high acidity and serious definicieny in nutrients (Lopes 1996). Improved
management practicies such as liming and fertilizerization have greatly increased soil fertility for growing
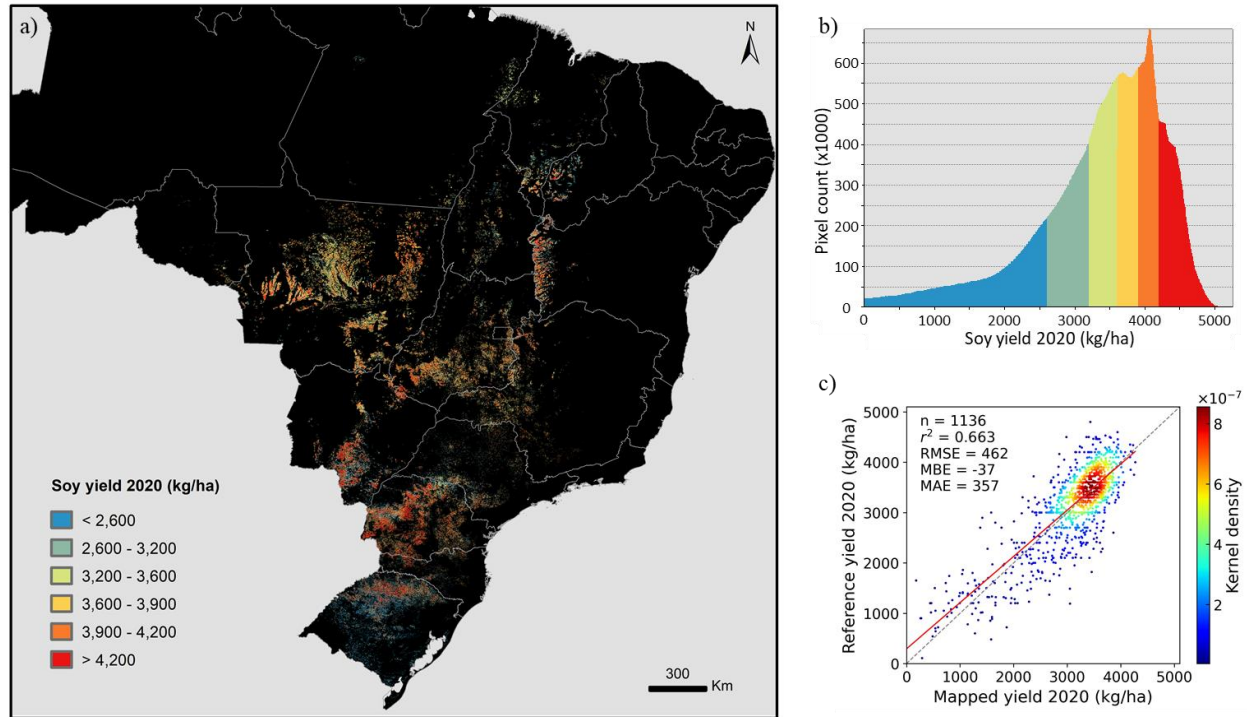
29

488  soybeans (Lopes 1996). These important changes in soil property are not reflected by the HWSD soil

489  database — likely the principal reason why the soil data did not contribute to soybean yield modelilng.

490  Crop modeling studies suggest that soil-related yield variability outweighs the simulated year-to-year

491  variations in yield due to weather when no fertilizer is applied (Folberth et al., 2016). Up-to-date high-

492  quality soil data may improve modeling yiled for soybean and other crops in the tropics where agriculture

493  is expanding (Eigenbrod et al. 2020). Future studies will investigate the utility of higher resolution soil

494  dataset for yield mapping (Hengl et al. 2017). Generating other spatially explicit data on agricultural

495  management that are important for crop production such as seed variety and fertilizer use, is another

496  potential way of improving yield mapping.

497  Lastly, a common practice in crop yield mapping is to construct a machine learning model at an

498  aggregated spatial scale where public yield statsitics are available, and apply the model to a finer scale at

499  which remote sensing data are acquired (e.g. Johnson 2014). The upscaling process (e.g. spatial

500  aggregation from pixel to municipal) can reduce uncertainties in the original data, as pixel-level errors

501  may be averaged out. Our yield models were calibrated at the municipal scale. More problematic is the

502  downscaling process (i.e. applying the trained model to pixels), as pixel-level errors often exist from e.g.

503  atmospheric correction or misclassification. The discrepency between model performance (Figure 5,

504  overall RMSE 344 kg/ha) and yield map assessment at the same municipality scale (Figure 10, overall

505  RMSE 418 kg/ha) revealed a positive bias in the predicted yield (Figure 9), although the models were

506  unbiased after linear adjustment (Figure 4). This bias was primarily stemmed from the downscaling

507  process, where pixel-level errors couls corrupt the results. Such bias may be removed using field-based

508  yield measurements. However, such datasets are traditionally held by private industry without public

509  access especially over large areas such as the national scale (see Deines et al. (2021) for the case of the

510  United States). Open access to field observations is rare in most parts of the world (Coutu et al. 2020).

511  Increasing the access to historical field observations is a potentially effective way of advancing crop yield

512  research.

513    **4.2.Towards operational yield mapping**

514    Achieving operational yield prediction using satellite data alone is a cost-effective approach of generating

515    timely information on crop production. To demonstrate the predictive capability of our yield models, we

516    applied the models, trained on 2001-2019 data, to 2020 data and produced a 30 m resolution soybean

517    yield map for 2020 (Figure 14). We also collected municipal yield statistics for 2020 and compared with

518    our 2020 yield map. Our random forests models, trained on 2001-2019 data, were able to predict 2020

519    yield with comparable accuracy as the withheld 2001-2019 test data. The RMSE, MBE and $r^2$ of the

520    direct output of random forests predictions for 2020 was 555 kg/ha, -145 kg/ha and 0.66, respectively.

521    Consistent with the model performance on 2001-2019 test data, an overall bias was noted. To eliminate

522    this bias, we applied the linear regression approach as reported above. We randomly selected 3% of

523    municipalities (n=34) from the 1,136 municipalities, and constructed a linear regression model using the

524    random forests-predicted yield as the independent variable and the 2020 municipal yield statistics as the

525    dependent variable. After bias correction, the MBE was reduced to -37 kg/ha, and RMSE was reduced to

526    462 kg/ha (Figure 14b). The RMSE represents 13% error relative to the national average of 3480 kg/ha in

527    2020. This result suggests that our pre-trained models can be used to generate high-resolution soybean

528    yield maps for future years with the caveat that a small amount of reference data are still needed for the

529    final bias correction. Given the continued operational satellite data acquisitions, including Landsat 8,

530    Landsat 9, MODIS and Visible Infrared Imaging Radiometer Suite (VIIRS), the demonstrated predictive

531    capability of our pre-trained yield models may be used for future yield mapping in a semi-operational

532    mode.

**Figure 14.** Soybean yield in year 2020 predicted using models trained on 2001-2019 data. a) 30-m map of soy yield 2020. b) Density distribution of the soy yield map. The colors match those shown on the map, and each color corresponds to approximately 1/6 of the total soy pixels. c) Comparison between predicted yield and municipal yield statistics as reference.

The rapidly developing technology of satellite remote sensing is transforming global agriculture. Earth observation data are increasingly used in research and operational settings for mapping crop types, monitoring crop growth, improving agricultural management and forecasting food production. Increasing the comprehensiveness within a single data product, including area, yield, cropping intensity and calendar, at high spatial and temporal resolution has been identified as one of the future research areas in developing global gridded cropping system data product (Kim et al. 2021). We showed in a previous study that satellite data could be used retrospectively mapping soybean over South America since 2001 (Song et al. 2021a). Our 30 m South America soybean map product is being updated at an annual frequency in an operatioanl mode as new satellite data are acquired. This study extends our research from

548  crop type mapping to yield mapping, and we demonstrated that pre-trained machine learning models

549  could be applied for yield mapping in future years. Our current approach for yield mapping and updating

550  uses satellite data of the entire growing season as input. This post-season mapping can generate highly

551  relabile data products, but lacks sufficient timeliness to capture production shocks resulted from e.g.

552  extreme weather events within the growing season. Recent research has demonstrated that early- and in-

553  season crop type mapping and crop yield forecasting could be achieved using advanced machine learning

554  algorithms (e.g. Lin et al. 2022), seasonal climate forecast (Iizumi et al. 2021), and in-season weather

555  observations (Schauberger et al. 2017). Implementing robust in-season forecasting methods in monitoring

556  systems is needed to mitigate the adverse impacts of climate change (Fritz et al. 2019; Kim et al. 2021; Li

557  et al. 2019; Lobell and Burke 2010; Nakalembe et al. 2021).


558  **5. Conclusions**

559  We developed a machine learning-based approach to map annual soybean yield in Brazil over the past

560  two decades. Consistent satellite observations from the open Landsat and MODIS data archives were used

561  to calibrate unbiased yield models using random forests followed by linear regression. Soybean yield

562  maps were generated at 30-meter spatial resolution for every year from 2001 to 2020. NDVI at the peak

563  of the growing season was found to be the most important variable for modeling soybean yield. Our study

564  explicitly demonstrated the utility of climate and weather variables for crop yield estimation. Our multi-

565  scale approach was effective in integrating official yield statistics at political unit level with remote

566  sensing data. Our study demonstrated that models trained on long-term historical data could be employed

567  to predict yield for future years. Our research also highlights that improving the temporal density of high-

568  resolution satellite observations, and enhancing the accessibility to field-level yield measurements are

569  viable ways to improve crop yield mapping over large areas.

570

571

576 **References**

577 Battude, M., Al Bitar, A., Morin, D., Cros, J., Huc, M., Marais Sicre, C., Le Dantec, V., & Demarez, V.

578 (2016). Estimating maize biomass and yield over large areas using high spatial and temporal

579 resolution Sentinel-2 like remote sensing data. *Remote Sensing of Environment, 184*, 668-681

580 Becker-Reshef, I., Vermote, E., Lindeman, M., & Justice, C. (2010). A generalized regression-based

581 model for forecasting winter wheat yields in Kansas and Ukraine using MODIS data. *Remote*

582 *Sensing of Environment, 114*, 1312-1323

583 Benami, E., Jin, Z., Carter, M.R., Ghosh, A., Hijmans, R.J., Hobbs, A., Kenduiywo, B., & Lobell, D.B.

584 (2021). Uniting remote sensing, crop modelling and economics for agricultural risk management.

585 *Nature Reviews Earth & Environment, 2*, 140-159

586 Bolton, D.K., & Friedl, M.A. (2013). Forecasting crop yield using remotely sensed vegetation indices and

587 crop phenology metrics. *Agricultural and Forest Meteorology, 173*, 74-84

588 Breiman, L. (2001). Random Forests. *Machine Learning, 45*, 5-32

589 Cai, Y., Guan, K., Lobell, D., Potgieter, A.B., Wang, S., Peng, J., Xu, T., Asseng, S., Zhang, Y., You, L.,

590 & Peng, B. (2019). Integrating satellite and climate data to predict wheat yield in Australia using

591 machine learning approaches. *Agricultural and Forest Meteorology, 274*, 144-159

592 Claverie, M., Demarez, V., Duchemin, B., Hagolle, O., Ducrot, D., Marais-Sicre, C., Dejoux, J.-F., Huc,

593 M., Keravec, P., Béziat, P., Fieuzal, R., Ceschia, E., & Dedieu, G. (2012). Maize and sunflower

594 biomass estimation in southwest France using high spatial and temporal resolution remote sensing

595 data. *Remote Sensing of Environment, 124*, 844-857

596    Claverie, M., Ju, J., Masek, J.G., Dungan, J.L., Vermote, E.F., Roger, J.-C., Skakun, S.V., & Justice, C.

597         (2018). The Harmonized Landsat and Sentinel-2 surface reflectance data set. *Remote Sensing of*

598         *Environment, 219*, 145-161

599    Coutu, S., Becker-Reshef, I., Whitcraft, A.K., & Justice, C. (2020). Food security: underpin with public

600         and private data sharing. *Nature, 578*, 515

601    de Wit, A., Boogaard, H., Fumagalli, D., Janssen, S., Knapen, R., van Kraalingen, D., Supit, I., van der

602         Wijngaart, R., & van Diepen, K. (2019). 25 years of the WOFOST cropping systems model.

603         *Agricultural Systems, 168*, 154-167

604    de Wit, A., Duveiller, G., & Defourny, P. (2012). Estimating regional winter wheat yield with WOFOST

605         through the assimilation of green area index retrieved from MODIS observations. *Agricultural*

606         *and Forest Meteorology, 164*, 39-52

607    Defourny, P., Bontemps, S., Bellemans, N., Cara, C., Dedieu, G., Guzzonato, E., Hagolle, O., Inglada, J.,

608         Nicola, L., Rabaute, T., Savinaud, M., Udroiu, C., Valero, S., Bégué, A., Dejoux, J.-F., El Harti,

609         A., Ezzahar, J., Kussul, N., Labbassi, K., Lebourgeois, V., Miao, Z., Newby, T., Nyamugama, A.,

610         Salh, N., Shelestov, A., Simonneaux, V., Traore, P.S., Traore, S.S., & Koetz, B. (2019). Near

611         real-time agriculture monitoring at national scale at parcel resolution: Performance assessment of

612         the Sen2-Agri automated system in various cropping systems around the world. *Remote Sensing*

613         *of Environment, 221*, 551-568

614    DeFries, R.S., Field, C.B., Fung, I., Justice, C.O., Los, S., Matson, P.A., Matthews, E., Mooney, H.A.,

615         Potter, C.S., Prentice, K., Sellers, P.J., Townshend, J.R.G., Tucker, C.J., Ustin, S.L., & Vitousek,

616         P.M. (1995). Mapping the land surface for global atmosphere‑biosphere models: Toward

617         continuous distributions of vegetation's functional properties. *Journal of Geophysical Research:*

618         *Atmospheres, 100,* 20867-20882.

619  Deines, J.M., Patel, R., Liang, S.-Z., Dado, W., & Lobell, D.B. (2021). A million kernels of truth: Insights

620      into scalable satellite maize yield mapping and yield gap analysis from an extensive ground

621      dataset in the US Corn Belt. *Remote Sensing of Environment, 253*, 112174

622  Delécolle, R., Maas, S., Guerif, M., & Baret, F. (1992). Remote sensing and crop production models:

623      present trends. *ISPRS Journal of Photogrammetry and Remote Sensing, 47*, 145-161

624  Doraiswamy, P.C., Hatfieldb, J.L., Jacksona, T.J., Akhmedova, B., Pruegerb, J., & Sterna, A. (2004).

625      Crop condition and yield simulations using Landsat and MODIS. *Remote Sensing of*

626      *Environment, 92*, 548-559

627  Duchemin, B., Maisongrande, P., Boulet, G., & Benhadj, I. (2008). A simple algorithm for yield

628      estimates: Evaluation for semi-arid irrigated winter wheat monitored with green leaf area index.

629      *Environmental Modelling & Software, 23*, 876-892

630  Eigenbrod, F., Beckmann, M., Dunnett, S., Graham, L., Holland, R.A., Meyfroidt, P., Seppelt, R., Song,

631      X.P., Spake, R., Vaclavik, T., & Verburg, P.H. (2020). Identifying Agricultural Frontiers for

632      Modeling Global Cropland Expansion. *One Earth, 3*, 504-514

633  FAO (2020). FAOSTAT database. In. Rome: FAO

634  FAO/IIASA/ISRIC/ISSCAS/JRC (2012). Harmonized World Soil Database (version 1.2). In. FAO,

635      Rome, Italy and IIASA, Laxenburg, Austria.

636  Folberth, C., Skalsky, R., Moltchanova, E., Balkovic, J., Azevedo, L.B., Obersteiner, M., & van der

637      Velde, M. (2016). Uncertainty in soil data can outweigh climate impact signals in global crop

638      yield simulations. *Nat Commun, 7*, 11872

639  Franch, B., Vermote, E.F., Becker-Reshef, I., Claverie, M., Huang, J., Zhang, J., Justice, C., & Sobrino,

640      J.A. (2015). Improving the timeliness of winter wheat production forecast in the United States of

641      America, Ukraine and China using MODIS data and NCAR Growing Degree Day information.

642      *Remote Sensing of Environment, 161*, 131-148

643  Fritz, S., See, L., Bayas, J.C.L., Waldner, F., Jacques, D., Becker-Reshef, I., Whitcraft, A., Baruth, B.,

644      Bonifacio, R., Crutchfield, J., Rembold, F., Rojas, O., Schucknecht, A., Van der Velde, M.,

645        Verdin, J., Wu, B., Yan, N., You, L., Gilliams, S., Mücher, S., Tetrault, R., Moorthy, I., &

646        McCallum, I. (2019). A comparison of global agricultural monitoring systems and current gaps.

647        *Agricultural Systems, 168*, 258-272

648  Fritz, S., See, L., McCallum, I., You, L., Bun, A., Moltchanova, E., Duerauer, M., Albrecht, F., Schill, C.,

649        Perger, C., Havlik, P., Mosnier, A., Thornton, P., Wood-Sichra, U., Herrero, M., Becker-Reshef,

650        I., Justice, C., Hansen, M., Gong, P., Abdel Aziz, S., Cipriani, A., Cumani, R., Cecchi, G.,

651        Conchedda, G., Ferreira, S., Gomez, A., Haffani, M., Kayitakire, F., Malanding, J., Mueller, R.,

652        Newby, T., Nonguierma, A., Olusegun, A., Ortner, S., Rajak, D.R., Rocha, J., Schepaschenko, D.,

653        Schepaschenko, M., Terekhov, A., Tiangwa, A., Vancutsem, C., Vintrou, E., Wenbin, W., van

654        der Velde, M., Dunwoody, A., Kraxner, F., & Obersteiner, M. (2015). Mapping global cropland

655        and field size. *Glob Chang Biol, 21*, 1980-1992

656  Funk, C., & Budde, M.E. (2009). Phenologically-tuned MODIS NDVI-based production anomaly

657        estimates for Zimbabwe. *Remote Sensing of Environment, 113*, 115-125

658  Funk, C., Peterson, P., Landsfeld, M., Pedreros, D., Verdin, J., Shukla, S., Husak, G., Rowland, J.,

659        Harrison, L., Hoell, A., & Michaelsen, J. (2015). The climate hazards infrared precipitation with

660        stations--a new environmental record for monitoring extremes. *Sci Data, 2*, 150066

661  Gallego, F.J. (2004). Remote sensing and land cover area estimation. *International Journal of Remote*

662        *Sensing, 25*, 3019-3047

663  Gao, F., Anderson, M., Daughtry, C., & Johnson, D. (2018). Assessing the Variability of Corn and

664        Soybean Yields in Central Iowa Using High Spatiotemporal Resolution Multi-Satellite Imagery.

665        *Remote Sensing*, *10,* 1489.

666  Gibbs, H.K., Ruesch, A.S., Achard, F., Clayton, M.K., Holmgren, P., Ramankutty, N., & Foley, J.A.

667        (2010). Tropical forests were the primary sources of new agricultural land in the 1980s and

668        1990s. *Proceedings of the National Academy of Sciences, USA, 107*, 16732-16737

669  Gitelson, A.A. (2004). Wide Dynamic Range Vegetation Index for remote quantification of biophysical

670        characteristics of vegetation. *Journal of Plant Physiology, 161*, 165-173

671    Gonzáles-Alonso, F., & Cuevas, J.M. (1993). Remote sensing and agricultural statistics: crop area

672        estimation through regression estimators and confusion matrices. *International Journal of Remote*

673        *Sensing, 14*, 1215-1219

674    Gorelick, N., Hancher, M., Dixon, M., Ilyushchenko, S., Thau, D., & Moore, R. (2017). Google Earth

675        Engine: Planetary-scale geospatial analysis for everyone. *Remote Sensing of Environment, 202*,

676        18-27

677    Hansen, M.C., Potapov, P.V., Moore, R., Hancher, M., Turubanova, S.A., Tyukavina, A., Thau, D.,

678        Stehman, S.V., Goetz, S.J., Loveland, T.R., & Kommareddy, A. (2013). High-resolution global

679        maps of 21st-century forest cover change. *Science, 342,* 850-853.

680    Harris, I., Osborn, T.J., Jones, P., & Lister, D. (2020). Version 4 of the CRU TS monthly high-resolution

681        gridded multivariate climate dataset. *Sci Data, 7*, 109

682    Hengl, T., Mendes de Jesus, J., Heuvelink, G.B., Ruiperez Gonzalez, M., Kilibarda, M., Blagotic, A.,

683        Shangguan, W., Wright, M.N., Geng, X., Bauer-Marschallinger, B., Guevara, M.A., Vargas, R.,

684        MacMillan, R.A., Batjes, N.H., Leenaars, J.G., Ribeiro, E., Wheeler, I., Mantel, S., & Kempen,

685        B. (2017). SoilGrids250m: Global gridded soil information based on machine learning. *PLoS*

686        *One, 12*, e0169748

687    Holzworth, D.P., Huth, N.I., deVoil, P.G., Zurcher, E.J., Herrmann, N.I., McLean, G., Chenu, K., van

688        Oosterom, E.J., Snow, V., Murphy, C., Moore, A.D., Brown, H., Whish, J.P.M., Verrall, S.,

689        Fainges, J., Bell, L.W., Peake, A.S., Poulton, P.L., Hochman, Z., Thorburn, P.J., Gaydon, D.S.,

690        Dalgliesh, N.P., Rodriguez, D., Cox, H., Chapman, S., Doherty, A., Teixeira, E., Sharp, J.,

691        Cichota, R., Vogeler, I., Li, F.Y., Wang, E., Hammer, G.L., Robertson, M.J., Dimes, J.P.,

692        Whitbread, A.M., Hunt, J., van Rees, H., McClelland, T., Carberry, P.S., Hargreaves, J.N.G.,

693        MacLeod, N., McDonald, C., Harsdorf, J., Wedgwood, S., & Keating, B.A. (2014). APSIM –

694        Evolution towards a new generation of agricultural systems simulation. *Environmental Modelling*

695        *& Software, 62*, 327-350

696    Hosseini, M., Kerner, H.R., Sahajpal, R., Puricelli, E., Lu, Y.-H., Lawal, A.F., Humber, M.L., Mitkish,

697         M., Meyer, S., & Becker-Reshef, I. (2020). Evaluating the Impact of the 2020 Iowa Derecho on

698         Corn and Soybean Fields Using Synthetic Aperture Radar. *Remote Sensing, 12*, 3878

699    Hu, Q., Yin, H., Friedl, M.A., You, L., Li, Z., Tang, H., & Wu, W. (2021). Integrating coarse-resolution

700         images and agricultural statistics to generate sub-pixel crop type maps and reconciled area

701         estimates. *Remote Sensing of Environment, 258*, 112365

702    Huang, J., Tian, L., Liang, S., Ma, H., Becker-Reshef, I., Huang, Y., Su, W., Zhang, X., Zhu, D., & Wu,

703         W. (2015). Improving winter wheat yield estimation by assimilation of the leaf area index from

704         Landsat TM and MODIS data into the WOFOST model. *Agricultural and Forest Meteorology,*

705         *204*, 106-121

706    Huang, X., Schneider, A., & Friedl, M.A. (2016). Mapping sub-pixel urban expansion in China using

707         MODIS and DMSP/OLS nighttime lights. *Remote Sensing of Environment, 175,* 92-108

708    Iizumi, T., Shin, Y., Choi, J., van der Velde, M., Nisini, L., Kim, W., & Kim, K.-H. (2021). Evaluating

709         the 2019 NARO-APCC Joint Crop Forecasting Service Yield Forecasts for Northern Hemisphere

710         Countries. *Weather and Forecasting, 36*, 879-891

711    Ines, A.V.M., Das, N.N., Hansen, J.W., & Njoku, E.G. (2013). Assimilation of remotely sensed soil

712         moisture and vegetation with a crop simulation model for maize yield prediction. *Remote Sensing*

713         *of Environment, 138*, 149-164

714    Ishikawa, Y., & Yamazaki, D. (2021). Global high-resolution estimation of cropland suitability and its

715         comparative analysis to actual cropland distribution. *Hydrological Research Letters, 15*, 9-15

716    Jin, X., Kumar, L., Li, Z., Feng, H., Xu, X., Yang, G., & Wang, J. (2018). A review of data assimilation

717         of remote sensing and crop models. *European Journal of Agronomy, 92*, 141-152

718    Jin, Z., Azzari, G., You, C., Di Tommaso, S., Aston, S., Burke, M., & Lobell, D.B. (2019). Smallholder

719         maize area and yield mapping at national scales with Google Earth Engine. *Remote Sensing of*

720         *Environment, 228*, 115-128

721 Johnson, D.M. (2014). An assessment of pre- and within-season remotely sensed variables for forecasting

722     corn and soybean yields in the United States. *Remote Sensing of Environment, 141*, 116-128

723 Johnson, M.D., Hsieh, W.W., Cannon, A.J., Davidson, A., & Bédard, F. (2016). Crop yield forecasting on

724     the Canadian Prairies by remotely sensed vegetation indices and machine learning methods.

725     *Agricultural and Forest Meteorology, 218-219*, 74-84

726 Jones, J.W., Hoogenboom, G., Porter, C.H., Boote, K.J., Batchelor, W.D., Hunt, L.A., Wilkens, P.W.,

727     Singh, U., Gijsman, A.J., & Ritchie, J.T. (2003). The DSSAT cropping system model. *European*

728     *Journal of Agronomy, 18*, 235–265

729 Kang, Y., & Özdoğan, M. (2019). Field-level crop yield mapping with Landsat using a hierarchical data

730     assimilation approach. *Remote Sensing of Environment, 228*, 144-163

731 Kim, K.-H., Doi, Y., Ramankutty, N., & Iizumi, T. (2021). A review of global gridded cropping system

732     data products. *Environmental Research Letters, 16*

733 King, L., Adusei, B., Stehman, S., Potapov, P.V., Song, X.-P., Krylov, A., Bella, C.D., Loveland, T.R.,

734     Johnson, D.M., & Hansen, M.C. (2017). A multi-resolution approach to national-scale cultivated

735     area estimation of soybean. *Remote Sensing of Environment, 195*, 13-29

736 Li, A., Liang, S., Wang, A., & Qin, J. (2007). Estimating crop yield from multi-temporal satellite data

737     using multivariate regression and neural network techniques. *Photogrammetric Engineering &*

738     *Remote Sensing, 73*, 1149-1157

739 Li, X.-Y., Li, X., Fan, Z., Mi, L., Kandakji, T., Song, Z., Li, D., & Song, X.-P. (2022). Civil war hinders

740     crop production and threatens food security in Syria. *Nature Food, 3*, 38-46

741 Li, Y., Guan, K., Schnitkey, G.D., DeLucia, E., & Peng, B. (2019). Excessive rainfall leads to maize yield

742     loss of a comparable magnitude to extreme drought in the United States. *Glob Chang Biol, 25*,

743     2325-2337

744 Liu, J., Huffman, T., Qian, B., Shang, J., Li, Q., Dong, T., Davidson, A., & Jing, Q. (2020). Crop yield

745     estimation in the Canadian Prairies using Terra/MODIS-derived crop metrics. *IEEE Journal of*

746     *Selected Topics in Applied Earth Observations and Remote Sensing, 13,* 2685-2697

747     Lobell, D.B. (2013). The use of satellite data for crop yield gap analysis. *Field Crops Research, 143*, 56-

748         64

749     Lobell, D.B., & Burke, M.B. (2010). On the use of statistical models to predict crop yield responses to

750         climate change. *Agricultural and Forest Meteorology, 150*, 1443-1452

751     Lobell, D.B., Thau, D., Seifert, C., Engle, E., & Little, B. (2015). A scalable satellite-based crop yield

752         mapper. *Remote Sensing of Environment, 164*, 324-333

753     Lopes, A.S. (1996). Soils under Cerrado: a success story in soil management. *Better crops international,*

754         *10*, 10

755     Massey, R., Sankey, T.T., Congalton, R.G., Yadav, K., Thenkabail, P.S., Ozdogan, M., & Sánchez

756         Meador, A.J. (2017). MODIS phenology-derived, multi-year distribution of conterminous U.S.

757         crop types. *Remote Sensing of Environment, 198*, 490-503

758     Mateo-Sanchis, A., Piles, M., Munoz-Mari, J., Adsuara, J.E., Perez-Suay, A., & Camps-Valls, G. (2019).

759         Synergistic integration of optical and microwave satellite data for crop yield estimation. *Remote*

760         *Sens Environ, 234*, 111460

761     Moulin, S., Bondeau, A., & Delecolle, R. (1998). Combining agricultural crop models and satellite

762         observations: From field to regional scales. *International Journal of Remote Sensing, 19*, 1021-

763         1036

764     Mulla, D.J. (2013). Twenty five years of remote sensing in precision agriculture: Key advances and

765         remaining knowledge gaps. *Biosystems Engineering, 114*, 358-371

766     Nakalembe, C., Becker-Reshef, I., Bonifacio, R., Hu, G., Humber, M.L., Justice, C.J., Keniston, J.,

767         Mwangi, K., Rembold, F., Shukla, S., Urbano, F., Whitcraft, A.K., Li, Y., Zappacosta, M., Jarvis,

768         I., & Sanchez, A. (2021). A review of satellite-based global agricultural monitoring systems

769         available for Africa. *Global Food Security, 29*

770     Nearing, G.S., Crow, W.T., Thorp, K.R., Moran, M.S., Reichle, R.H., & Gupta, H.V. (2012).

771         Assimilating remote sensing observations of leaf area index and soil moisture for wheat yield

772         estimates: An observing system simulation experiment. *Water Resources Research, 48*, W05525

773   Paudel, D., Boogaard, H., de Wit, A., Janssen, S., Osinga, S., Pylianidis, C., & Athanasiadis, I.N. (2021).

774         Machine learning for large-scale crop yield forecasting. *Agricultural Systems, 187*, 103016

775   Potapov, P., Hansen, M.C., Kommareddy, I., Kommareddy, A., Turubanova, S., Pickens, A., Adusei, B.,

776         Tyukavina, A., & Ying, Q. (2020). Landsat Analysis Ready Data for Global Land Cover and

777         Land Cover Change Mapping. *Remote Sensing, 12*, 426

778   Potapov, P., Turubanova, S., Hansen, M.C., Tyukavina, A., Zalles, V., Khan, A., Song, X.-P., Pickens,

779         A., Shen, Q., & Cortez, J. (2022). Global maps of cropland extent and change show accelerated

780         cropland expansion in the twenty-first century. *Nature Food, 3*, 19-28

781   Sakamoto, T., Gitelson, A.A., & Arkebauer, T.J. (2013). MODIS-based corn grain yield estimation model

782         incorporating crop phenology information. *Remote Sensing of Environment, 131*, 215-231

783   Schauberger, B., Gornott, C., & Wechsung, F. (2017). Global evaluation of a semiempirical model for

784         yield anomalies and application to within-season yield forecasting. *Glob Chang Biol, 23*, 4750-

785         4764

786   Schwalbert, R.A., Amado, T., Corassa, G., Pott, L.P., Prasad, P.V.V., & Ciampitti, I.A. (2020). Satellite-

787         based soybean yield forecast: Integrating machine learning and weather data for improving crop

788         yield prediction in southern Brazil. *Agricultural and Forest Meteorology, 284*, 107886

789   Shahhosseini, M., Hu, G., Huber, I., & Archontoulis, S.V. (2021). Coupling machine learning and crop

790         modeling improves crop yield prediction in the US Corn Belt. *Sci Rep, 11*, 1606

791   Skakun, S., Franch, B., Vermote, E., Roger, J.-C., Becker-Reshef, I., Justice, C., & Kussul, N. (2017).

792         Early season large-area winter crop mapping using MODIS NDVI data, growing degree days

793         information and a Gaussian mixture model. *Remote Sensing of Environment, 195*, 244-258

794   Skakun, S., Kalecinski, N.I., Brown, M.G.L., Johnson, D.M., Vermote, E.F., Roger, J.-C., & Franch, B.

795         (2021). Assessing within-Field Corn and Soybean Yield Variability from WorldView-3, Planet,

796         Sentinel-2, and Landsat 8 Satellite Imagery. *Remote Sensing, 13*, 872

797   Song, X.-P., Hansen, M.C., Potapov, P.V., Adusei, B., Pickering, J., Adami, M., Lima, A., Zalles, V.,

798         Stehman, S.V., Di Bella, C.M., Conde, M.C., Copati, E.J., Fernandes, L.B., Hernandez-Serna, A.,

799   Jantz, S.M., Pickens, A.H., Turubanova, S., & Tyukavina, A. (2021a). Massive soybean

800     expansion in South America since 2000 and implications for conservation. *Nature Sustainability,*

801     *4,* 784-792

802  Song, X.-P., Hansen, M.C., Stehman, S.V., Potapov, P.V., Tyukavina, A., Vermote, E.F., & Townshend,

803     J.R. (2018). Global land change from 1982 to 2016. *Nature, 560*, 639-643

804  Song, X.-P., Huang, W., Hansen, M.C., & Potapov, P. (2021b). An evaluation of Landsat, Sentinel-2,

805     Sentinel-1 and MODIS data for crop type mapping. *Science of Remote Sensing, 3*, 100018

806  Song, X.-P., Potapov, P.V., Krylov, A., King, L., Di Bella, C.M., Hudson, A., Khan, A., Adusei, B.,

807     Stehman, S.V., & Hansen, M.C. (2017). National-scale soybean mapping and area estimation in

808     the United States using medium resolution satellite imagery and field survey. *Remote Sensing of*

809     *Environment, 190*, 383-395

810  Tucker, C.J., Holben, B., Elgin Jr, J., & McMurtrey III, J. (1980). Relationship of spectral data to grain

811     yield variation. *Photogrammetric Engineering and Remote Sensing, 46*, 657-666

812  Veloso, A., Mermoz, S., Bouvet, A., Le Toan, T., Planells, M., Dejoux, J.-F., & Ceschia, E. (2017).

813     Understanding the temporal behavior of crops using Sentinel-1 and Sentinel-2-like data for

814     agricultural applications. *Remote Sensing of Environment, 199*, 415-426

815  Vermote, E., & Wolfe, R. (2015). MOD09GQ MODIS/Terra Surface Reflectance Daily L2G Global

816     250m SIN Grid V006 [Data set]. NASA EOSDIS Land Processes DAAC. Accessed 2021-11-27

817     from https://doi.org/10.5067/MODIS/MOD09GQ.006

818  Vroege, W., Vrieling, A., & Finger, R. (2021). Satellite support to insure farmers against extreme

819     droughts. *Nature Food, 2*, 215-217

820  Wardlow, B.D., & Egbert, S.L. (2008). Large-area crop mapping using time-series MODIS 250 m NDVI

821     data: An assessment for the U.S. Central Great Plains. *Remote Sensing of Environment, 112*,

822     1096-1116

823  Weiss, M., Jacob, F., & Duveiller, G. (2020). Remote sensing for agricultural applications: A meta-

824     review. *Remote Sensing of Environment, 236*

825    Wieder, W.R., Boehnert, J., Bonan, G.B., & Langseth, M. (2014). Regridded Harmonized World Soil
826           Database v1.2. In: Oak Ridge National Laboratory Distributed Active Archive Center, Oak
827           Ridge, Tennessee, USA.

828    Williams, J., Jones, C., Kiniry, J., & Spanel, D.A. (1989). The EPIC crop growth model. *Transactions of*
829           *the ASAE, 32*, 497-0511

830    Yang, H.S., Dobermann, A., Lindquist, J.L., Walters, D.T., Arkebauer, T.J., & Cassman, K.G. (2004).
831           Hybrid-maize—a maize simulation model that combines two crop modeling approaches. *Field*
832           *Crops Research, 87*, 131-154

833    Zhang, G., & Lu, Y., (2012). Bias-corrected random forests in regression. *Journal of Applied Statistics,*
834           *39,* 151-160

835    Zalles, V., Hansen, M.C., Potapov, P.V., Parker, D., Stehman, S.V., Pickens, A.H., Parente, L.L.,
836           Ferreira, L.G., Song, X.-P., Hernandez-Serna, A., & Kommareddy, I. (2021). Rapid expansion of
837           human impact on natural land in South America since 1985. *Science Advances, 7*, eabg1620

838    Zalles, V., Hansen, M.C., Potapov, P.V., Stehman, S.V., Tyukavina, A., Pickens, A., Song, X.-P., Adusei,
839           B., Okpa, C., Aguilar, R., John, N., & Chavez, S. (2019). Near doubling of Brazil's intensive row
840           crop area since 2000. *Proceedings of the National Academy of Sciences, USA, 116*, 428-435

841

842