# Assessing the benefit of satellite-based Solar-Induced Chlorophyll Fluorescence in crop yield prediction

Bin Peng[a,b,*], Kaiyu Guan[a,b,c,*], Wang Zhou[a], Chongya Jiang[a,c], Christian Frankenberg[d,e], Ying Sun[f], Liyin He[d], Philipp Köhler[d]

[a] Department of Natural Resources and Environmental Sciences, College of Agricultural, Consumer and Environmental Sciences, University of Illinois at Urbana-Champaign, Urbana, Illinois, USA
[b] National Center of Supercomputing Applications, University of Illinois at Urbana-Champaign, Urbana, Illinois, USA
[c] Center for Advanced Bioenergy and Bioproducts Innovation, University of Illinois at Urbana-Champaign, Urbana, Illinois, USA
[d] Division of Geological and Planetary Sciences, California Institute of Technology, Pasadena, CA, USA
[e] Jet Propulsion Laboratory, California Institute of Technology, Pasadena, CA, USA
[f] Soil and Crop Sciences Section, School of Integrative Plant Science, Cornell University, Ithaca, NY, USA

## ARTICLE INFO

## ABSTRACT

Large-scale crop yield prediction is critical for early warning of food insecurity, agricultural supply chain management, and economic market. Satellite-based Solar-Induced Chlorophyll Fluorescence (SIF) products have revealed hot spots of photosynthesis over global croplands, such as in the U.S. Midwest. However, to what extent these satellite-based SIF products can enhance the performance of crop yield prediction when benchmarking against other existing satellite data remains unclear. Here we assessed the benefits of using three satellite-based SIF products in yield prediction for maize and soybean in the U.S. Midwest: gap-filled SIF from Orbiting Carbon Observatory 2 (OCO-2), new SIF retrievals from the TROPOspheric Monitoring Instrument (TROPOMI), and the coarse-resolution SIF retrievals from the Global Ozone Monitoring Experiment-2 (GOME-2). The yield prediction performances of using SIF data were benchmarked with those using satellite-based vegetation indices (VIs), including normalized difference vegetation index (NDVI), enhanced vegetation index (EVI), and near-infrared reflectance of vegetation (NIRv), and land surface temperature (LST). Five machine-learning algorithms were used to build yield prediction models with both remote-sensing-only and climate-remote-sensing-combined variables. We found that high-resolution SIF products from OCO-2 and TROPOMI outperformed coarse-resolution GOME-2 SIF product in crop yield prediction. Using high-resolution SIF products gave the best forward predictions for both maize and soybean yields in 2018, indicating the great potential of using satellite-based high-resolution SIF products for crop yield prediction. However, using currently available high-resolution SIF products did not guarantee consistently better yield prediction performances than using other satellite-based remote sensing variables in all the evaluated cases. The relative performances of using different remote sensing variables in yield prediction depended on crop types (maize or soybean), out-of-sample testing methods (five-fold-cross-validation or forward), and record length of training data. We also found that using NIRv could generally lead to better yield prediction performance than using NDVI, EVI, or LST, and using NIRv could achieve similar or even better yield prediction performance than using OCO-2 or TROPOMI SIF products. We concluded that satellite-based SIF products could be beneficial in crop yield prediction with more high-resolution and good-quality SIF products accumulated in the future.

## 1. Introduction

Crop yield forecasting at a regional to global scale is important for early warning of food insecurity, agricultural supply chain management, and economic market prediction (Everingham et al. 2002; Hansen and Indeje 2004; Isengildina-Massa et al. 2008). Generally, a crop yield forecasting system can be based on either physical or statistical models. Physical-model-based approach usually uses a crop model to dynamically simulate crop growth and yield formation processes (Brown et al. 2018; Jones et al. 2017; Jones et al. 2003; Peng

---

et al. 2018a; Peng et al. 2020; Rosenzweig et al. 2013; Shelia et al. 2019). However, due to the complexity and relatively lower performance of these physical models at large scales, statistical models are widely used in operational large-scale crop yield forecasting systems (Chipanshi et al. 2015; Li et al. 2019; Newlands et al. 2014; Peng et al. 2018b).

Statistical crop yield models are data-driven, and thus the type, volume, as well as quality of input data are among the key factors determining the model performance. Earlier studies developing statistical models for crop yield forecasting mainly rely on environmental factors as inputs, such as climate and soil condition (Legler et al. 1999; Phillips et al. 1999; Potgieter et al. 2002; Qian et al. 2009). Later, satellite data has been proved to be beneficial in operational crop yield forecasting systems. Using satellite data only or adding satellite data upon environmental information can generally lead to better yield estimation than traditional statistical crop yield models only using environmental factors (Li et al. 2019). The application of various remote sensing products across a diverse spectral range in crop yield estimation has been extensively explored (Guan et al. 2017), including (but not limited to) surface reflectance (You et al. 2017), vegetation indices (Bolton and Friedl 2013; Cai et al. 2019; Chipanshi et al. 2015; Johnson 2014; Lobell et al. 2015; Newlands et al. 2014; Peng et al. 2018b), land surface temperature (LST) (Cai et al. 2017; Johnson 2014; Li et al. 2019; You et al. 2017), fraction of photosynthetically-active radiation (fPAR) (Bastiaanssen and Ali 2003; Jiang et al. 2004), gross primary productivity (GPP) (He et al. 2018), evapotranspiration (Anderson et al. 2016; Yang et al. 2018), active microwave based backscattering and passive microwave based vegetation optical depth (Chaparro et al. 2018; Guan et al. 2017).

Satellite-based Solar-induced Chlorophyll Fluorescence (SIF) has recently demonstrated to be effective in capturing the spatial and temporal variabilities of terrestrial carbon uptake (Frankenberg et al. 2011; Guanter et al. 2014; Joiner et al. 2013; Parazoo et al. 2013; Shiga et al. 2018). Although agricultural areas are always hot spots on the satellite-based SIF maps during the peak growing season, few studies have directly explored the use of satellite SIF data in crop yield estimation. Guanter et al. (2014) and Guan et al. (2016) were the first to indirectly link SIF retrieval and crop yield. They first estimated GPP through linear scaling with or without accounting for stoichiometry and photosynthetic pathways, and then benchmarked with aggregated net primary productivity (NPP) estimated from production of all crops, instead of yield for individual crops mainly due to the coarse spatial resolution (0.5 degree) of the Global Ozone Monitoring Experiment-2 (GOME-2) gridded SIF products used in their studies. A recent work using SIF products from the SCanning Imaging Absorption spectroMeter for Atmospheric CHartographY (SCIAMACHY) and GOME-2 for wheat yield prediction in Australia showed SIF was no better than Enhanced Vegetation Index (EVI), largely due to the low spatial and temporal resolution of SCIAMACHY and GOME-2 SIF products and also low signal-to-noise ratio in these coarse-resolution SIF products (Cai et al. 2019). More recently, the Orbiting Carbon Observatory 2 (OCO-2) and the TROPOspheric Monitoring Instrument (TROPOMI) can provide SIF retrievals at much higher spatial resolutions (Frankenberg et al. 2014; Köhler et al. 2018), which opens the opportunity for large-scale crop yield estimation using satellite-based SIF. Although the original OCO-2 SIF product has its limitations in sparse sampling swath and long revisit cycle (Sun et al. 2018), data-driven gap-filling of OCO-2 SIF can provide spatial continuous and high resolution (0.05 degree) SIF products (Li and Xiao 2019; Yu et al. 2018; Zhang et al. 2018). However, there are still no studies directly using these satellite-based high-resolution SIF products for crop yield prediction, and further comparing the performances of using satellite-based high-resolution SIF products with those of using coarse-resolution SIF products and traditional vegetation indices or LST in crop yield prediction. Therefore, the benefits of using satellite-based high-resolution SIF products in operational crop yield estimation remains unclear. Besides the above advancements in

generating new high-resolution SIF products, the near-infrared reflectance of vegetation (NIRv) (Badgley et al. 2017), a new combination of red and near-infrared band reflectance from Moderate Resolution Imaging Spectroradiometer (MODIS), has been found to be well correlated with SIF and can lead to improved GPP estimates relative to the normalized difference vegetation index (NDVI) and fPAR, which indicates that NIRv can also be potentially applied for crop yield estimation. To our best knowledge, there is still no studies testing the performance of NIRv in crop yield prediction.

This study aims at assessing the potential application of satellite-based high-resolution SIF products from gap-filled OCO-2 and TROPOMI in estimating maize and soybean yield in the U.S. Midwest. The yield prediction performance using satellite-based high-resolution SIF will be benchmarked with those using coarse-resolution GOME-2 SIF, several MODIS-based vegetation indices (NDVI, EVI, and NIRv), and LST. By conducting this assessment and intercomparison study, we want to answer the following questions: (1) Does high-resolution SIF products perform better than coarse-resolution SIF products in crop yield prediction? (2) Does high-resolution SIF products perform better than vegetation indices and LST in crop yield prediction? Answer to these questions could guide the development of operational yield prediction system using multi-source remote sensing data.

## 2. Study Area and Data

### 2.1. Study Area

We focused on rainfed maize and soybean yield estimation over 12 states in the U.S. Midwest, including Illinois, Indiana, Iowa, Kansas, Michigan, Minnesota, Missouri, Nebraska, North Dakota, Ohio, South Dakota, and Wisconsin (Fig. 1). In 2018, harvested area over these 12 states accounted for 87% and 83% of U.S. total harvested area for maize and soybean, respectively, which corresponded to 89% and 84% of U.S. total maize and soybean productions, respectively. The rainfed maize and soybean harvested area in these 12 states accounted for 66% and 73% of U.S. total harvested area, which corresponded to 69% and 74% of U.S. total productions in 2018 for corn and soybean, respectively (Fig. 2).

### 2.2. Historical crop yield and acreage data from USDA NASS

We obtained county-level harvested yield and acreage data for both rainfed maize and soybean from the U.S. Department of Agriculture (USDA) National Agricultural Statistics Service (NASS) Quick Stats Database (quickstats.nass.usda.gov). For counties without any irrigated yield records, their yield was considered to be rainfed. For irrigated counties, we only included yield records that were explicitly reported as "nonirrigated" as rainfed yield.

### 2.3. Historical climate data

Historical climate data were obtained from the Parameter-elevation Relationships on Independent Slopes Model (PRISM), which has a 4-km spatial resolution (Daly et al. 2008). We used monthly mean temperature (Tair), precipitation (Prec), and vapor pressure deficit (VPD) as climate variables in yield prediction models. Tair and Prec are commonly used in building crop yield prediction models as they represent the basic meteorological condition over a region (Li et al. 2019; Lobell and Burke 2010; Lobell et al. 2015; Peng et al. 2018b). VPD has been found to be a dominant factor in indicating crop water stress over the U.S. Midwest (Lobell et al. 2014). Though VPD is highly correlated with Tair, adding VPD upon Tair and Prec still improved the forecasting performance of yield (Peng et al. 2018b). The original 4-km data was aggregated to the county level without differentiating crop types.
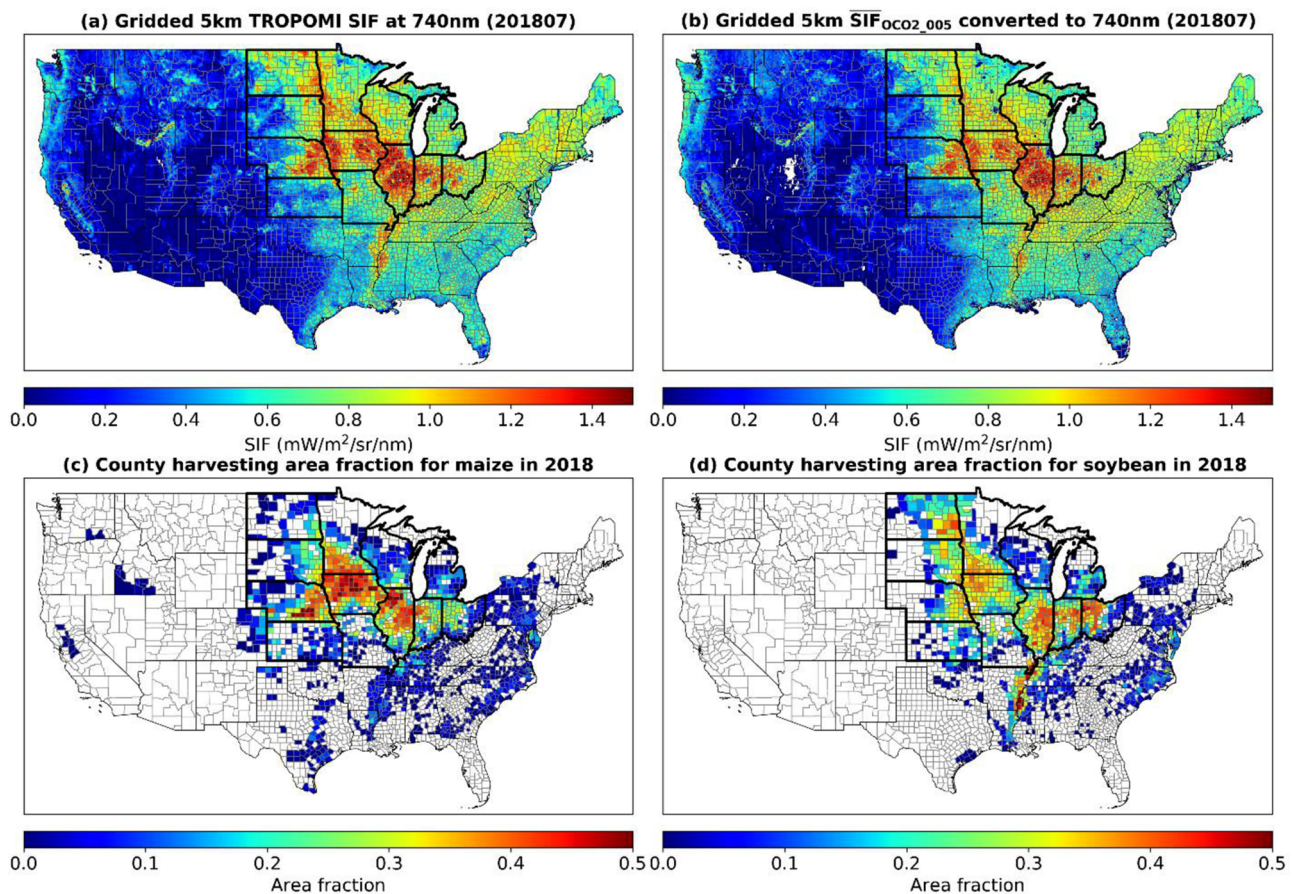
**Fig. 1.** Satellite-based monthly averaged SIF products reveal a summer photosynthesis hotspot over the U.S. Midwest, where maize and soybean fields dominate the landscape. (a) gridded 5 km TROPOMI SIF at 740 nm over the Contiguous United States (CONUS) in July, 2018; (b) gridded 5 km $\overline{SIF}_{OCO2\_005}$ converted to 740 nm over the CONUS in July, 2018; (c) county-level harvesting area fraction for maize in 2018; and (d) county-level harvesting area fraction for soybean in 2018. The 12 states in the U.S. Midwest are highlighted using thick black lines.

### 2.4. Satellite data

#### 2.4.1. Satellite-based SIF data

We used a spatially contiguous global OCO-2 SIF product at 0.05° and 16-day resolutions ($\overline{SIF}_{OCO2\_005}$) (Yu et al. 2018), which is a machine learning prediction of OCO-2 nadir SIF using MODIS Nadir Bidirectional Reflectance Distribution Function (BRDF)-Adjusted Reflectance (NBAR) products (MCD43A4 and MCD43C4). The original instantaneous OCO-2 SIF retrievals at 757 nm and 771 nm were scaled to 757 nm using ($SIF_{757}$ + 1.5 × $SIF_{771}$)/2 to improve the accuracy and finally converted to daily mean SIF before used for model training and prediction. Biome- and time-step-specific feedforward artificial neural network (ANN) models were trained and cross-validated using co-located OCO-2 footprints and MODIS NBAR data. This product shows high quality when benchmarked with independent airborne SIF measurements (Yu et al. 2018). The data is available after September of 2014, and we used all the data from 2015 to 2018.

We also used the TROPOMI ungridded daily SIF product with a footprint size of about 7 km x 3.5 km at nadir (Köhler et al. 2018). TROPOMI is on board the Sentinel 5 Precursor (S-5 P). SIF retrieval was conducted at a spectral window ranging from 743–758 nm, which is a subset of TROPOMI's band 6 (725–775 nm). TROPOMI SIF data is available after February of 2018. We gridded the footprint-level TROPOMI SIF data to 0.05° to match the spatial resolution of $\overline{SIF}_{OCO2\_005}$ product used in this study. A SIF value contributes to a grid cell average if the footprint covers the center of this grid cell (Köhler et al. 2018). We converted the instantaneous TROPOMI SIF to daily average by applying a day length correction factor contained in the TROPOMI SIF

product, which assumed cloud-free condition as a first-order approximation (Frankenberg et al. 2011; Köhler et al. 2018).

Besides $\overline{SIF}_{OCO2\_005}$ and TROPOMI SIF products, we also used a coarse-resolution SIF products from GOME-2 (Köhler et al. 2015). GOME-2 is on board EUMETSAT's polar orbiting Meteorological Operational Satellites (MetOp-A and MetOp-B), and a nadir-scanning medium-resolution UV/VIS spectrometer with a spectral range between 240 and 790 nm. A subchannel ranging from 720-758 nm were used to retrieve SIF signal at 740 nm from GOME-2 observations (Köhler et al. 2015). Before gridded into 0.5° product, the instantaneous SIF retrievals were converted to daily mean values using the daily correction factor approach (Frankenberg et al. 2011), which is the same with TROPOMI SIF product (Köhler et al. 2018). This product is available since 2007 and we used all data from 2015 to 2018.

Following Köhler et al. (2018), we scaled $\overline{SIF}_{OCO2\_005}$ at 757 nm to TROPOMI SIF retrieval channel (around 740 nm) by multiplying a factor of 1.56, which was determined based on a reference SIF emission shape derived from leaf-level measurements (Magney et al. 2017). The scaled $\overline{SIF}_{OCO2\_005}$, TROPOMI gridded SIF products at 740 nm, and GOME-2 SIF product were then aggregated to monthly and county level mean values for maize and soybean separately using crop fraction determined from the yearly cropland data layer (CDL, see section 2.4.3 for details).

#### 2.4.2. MODIS data

We used NDVI, EVI, NIRv, and LST data from MODIS as additional remote sensing based predictors for crop yield estimation. Both NDVI and EVI were from the Terra 16-day global vegetation indices product
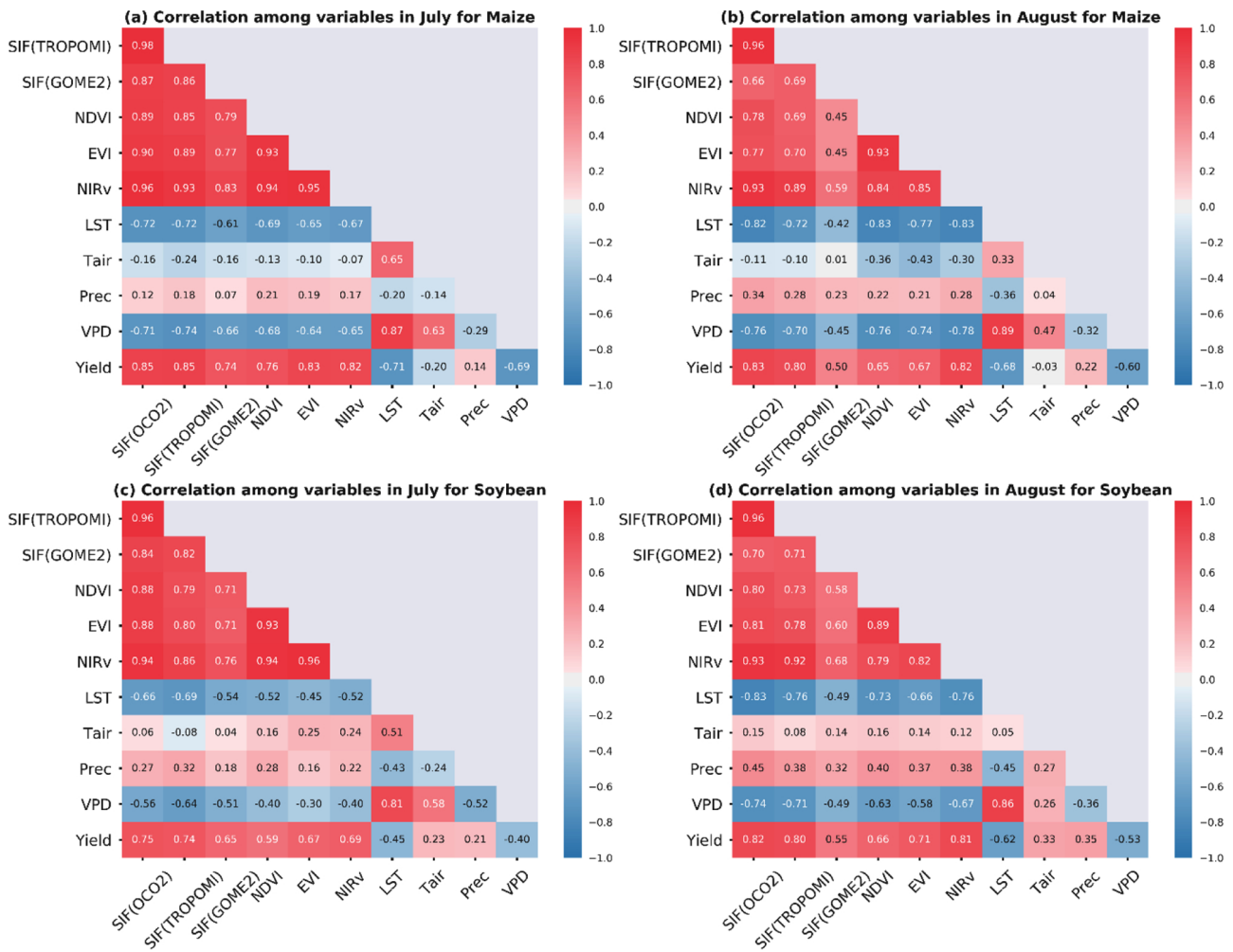
**Fig. 2.** Pearson correlation coefficients (r) among county-level climate and remote sensing variables in July (left column) and August (right column) as well as final harvested yield for maize (top) and soybean (bottom) in 2018. SIF(OCO2), SIF(TROPOMI), and SIF(GOME2) in the x-axis represent $\tilde{SIF}_{OCO2\_005}$, TROPOMI and GOME-2 SIF products, respectively.

with a spatial resolution of 250 m (MOD13Q1.006), which contains the best available vegetation index values from all the MODIS acquisitions within the 16-day period. NIRv was calculated using the daily MODIS Nadir BRDF-Adjusted Reflectance (NBAR) product with a spatial resolution of 500 m (MCD43A4.006) following the definition of NIRv = $(\rho_{NIR} - \rho_R)/(\rho_{NIR} + \rho_R) \times \rho_{NIR}$, where $\rho_{NIR}$ and $\rho_R$ represent the surface reflectance at near-infrared and red bands, respectively (Badgley et al. 2017). Daily NIRv was then composited into 16-day data following a similar maximum-value approach with NDVI and EVI. Daytime LST data was from the Aqua 8-day global land surface temperature and emissivity product with a spatial resolution of 1 km (MYD11A2.006). We choose Aqua daytime LST product here as the Aqua satellite goes cross the equator at approximately local time 1:30 P.M., which is closer to the time of maximum canopy temperature, incoming solar radiation, as well as most possible stressing conditions for crops on clear days, compared with Terra satellite with a visiting time of 10:30 A.M.. We did not use the nighttime LST product as it has little correlation with crop yield (Johnson 2014). All 8-day or 16-day MODIS data were firstly aggregated to monthly scale, and finally aggregated to county-level mean values for maize and soybean separately using crop area fractions determined from the yearly 30 m USDA NASS Cropland Data Layer (CDL, see section 2.4.3 for details). All MODIS pixels with area fraction of corn or soybean larger than 50% within a specific county were averaged as the county-level mean values.

### 2.4.3. CDL from USDA NASS

The USDA NASS CDL was used to aggregate the remote sensing variables to county level for corn and soybean separately. The CDL data is a yearly multi-satellite based crop type classification product using decision tree supervised classifier and has a 30 m spatial resolution. The classification accuracy for maize and soybean is above 95% over the U.S. Midwest (Boryan et al. 2011). For MODIS VIs and LST data, we aggregated all the MODIS pixels with fractions of corn or soybean larger than 50% within a county. For SIF data, we conducted simple weighted average of all the 5 km grids within a county using corn or soybean area fraction as weights.

## 3. Method

### 3.1. Crop yield model development

We used five different machine learning algorithms to develop the crop yield model, including the least absolute shrinkage and selection operator regression (LASSO) (Tibshirani 1996), ridge regression (RIDGE) (Hoerl and Kennard 1970), support vector regression (SVR) (Smola and Schölkopf 2004), random forest regression (RF) (Breiman 2001), and artificial neural network (ANN) (Gardner and Dorling 1998; Specht 1991). LASSO and RIDGE are both regularized regression methods and their difference is that LASSO uses L1 regularization while RIDGE uses L2 regularization (Fu 1998; Tibshirani 1996). Both LASSO

and RIDGE have the same penalty parameter $\alpha$ to be tuned. The SVR is a kernel-based regression method solving nonlinear regression problems by transferring the data to a higher-dimensional space through a kernel function. We used the radial basis function (RBF) kernel for SVR (Suykens and Vandewalle 1999) as it usually gives better accuracy than linear and polynomial kernels. RF is a binary-tree based machine learning algorithm, which builds an ensemble of decision trees with different subsets of variables. The ANN is based on a collection of artificial neurons, which loosely model the neurons in a biological brain and can receive inputs, change their internal states (activation) according to the inputs, and produce outputs depending on the inputs and activation. We used the multilayer perceptron (MLP) regressor (Gardner and Dorling 1998), which is feedforward ANN and trains using backpropagation with no activation function in the output layer. We choose the L-BFGS method to optimize the squared-loss as it converges faster and performs better for small datasets. All these methods have been previously explored for crop yield estimation at varied scales (Cai et al. 2019; Jeong et al. 2016; Jiang et al. 2004). The main purpose of using multiple algorithms with varied complexity here is to test whether the differences in yield predictability using different remote sensing variables are consistent when using different algorithms to build the crop yield model.

### 3.2. Experiment design

We conducted two groups of experiments: one group used only remote sensing variables, while another group used both climate and remote sensing variables as predictors. Remote sensing variables included monthly SIF, NDVI, EVI, NIRv, and LST during the growing season. We used monthly air temperature, precipitation, and VPD during the growing season as climate variables in the second group of experiments as these variables combined can provide reasonable prediction performance among all the climate-only models (Li et al. 2019; Peng et al. 2018b). The growing season in this study was defined as May to September, which aligned with the actual growing season of corn and soybean in the U.S. Midwest. All variables were standardized by removing their mean values and scaling to unit variance before used for model training and testing, which can help avoid bad performance if the individual features are not standard normally distributed data.

Two different out-of-sample validation methods were used to quantify the yield estimation performance. One was the repeated five-fold-cross-validation (FFCV) method, and the other one was the forward method. The repeated FFCV method runs the FFCV for $n$ times, each of which randomly splits the whole dataset into 5 folds, and uses 4 folds for training and 1 fold for testing within a FFCV loop. We choose $n = 100$ corresponding to 500 training-testing splits in total, which balanced well between accuracy and computation burden. The forward method used all data from years before the prediction year as training dataset. For both repeated FFCV and forward methods, all the five algorithms were automatically optimized by tuning their hyperparameters using FFCV on their training dataset. The training data was shuffled in a consistent way to avoid the impact of internal structures (both spatial and temporal) in training data on FFCV. The prediction performance was then assessed using the testing dataset. We used coefficient of determination ($R^2$), root mean square error (RMSE), and mean absolute bias (MAB) as statistical metrics in performance assessment. For repeated FFCV method, we reported both mean and standard deviation of these two metrics evaluated over the 500 training-testing splits.

To better demonstrate the benefit of using $S\bar{I}F_{OCO2\_005}$ and TROPOMI SIF in yield prediction, we evaluated the performance using data during 2015–2018 (4-year case hereafter) and only in 2018 (1-year case hereafter) mainly considering the data availability of $S\bar{I}F_{OCO2\_005}$ (2015–2018) and TROPOMI (2018 only) SIF products. For the 4-year case with FFCV validation methods, we only used $S\bar{I}F_{OCO2\_005}$ for model training and testing. For the 1-year case with FFCV validation method,

we trained and tested the models using $S\bar{I}F_{OCO2\_005}$ and TROPMI SIF in 2018 separately. For the forward method, the models were trained using $S\bar{I}F_{OCO2\_005}$ during 2015–2017 while tested using $S\bar{I}F_{OCO2\_005}$ and TROPOMI SIF in 2018 separately. For the sake of fair comparison, the data length for model training and testing using other remote sensing, climate, or combined variables was consistent with that of SIF products.

## 4. Results

### 4.1. Correlation of crop yield with climate and remote sensing variables in 2018

The spatial patterns of crop yield in 2018 were better correlated with remote sensing variables (SIF, NDVI, EVI, NIRv, and LST) than climate factors (Tair, Prec, and VPD). Among the tested remote sensing variables, SIF from OCO-2 and TROPOMI, EVI, NIRv in July and SIF from OCO-2 and TROPOMI, and NIRv in August showed correlation coefficients larger than 0.8 with maize yield. Similarly, SIF from OCO-2 and TROPOMI and NIRv in August also showed correlation coefficients larger than 0.8 with soybean yield. The correlation coefficient between SIF from GOME-2 and crop yield was smaller than those between SIF from OCO-2 and TROPOMI and crop yield, sometimes even ranked the lowest among all the remote sensing variables, such as in August for both maize and soybean. LST negatively correlated with crop yield. For maize, the correlation coefficient between LST and yield in July is larger than that in August. For soybean, the correlation coefficient between LST and yield in August is larger than that in July. Among the three climate factors, VPD was negatively correlated with crop yield and precipitation was positively correlated with crop yield for both corn and soybean, while air temperature was negatively correlated with crop yield for maize and positively correlated with crop yield for soybean. VPD showed higher correlation coefficients than Tair and Prec with both maize and soybean yields. For example, the correlation coefficients between VPD and maize yield were -0.69 and -0.60 in July and August, respectively; while those between Tair and maize yield were only -0.20 and -0.03 in July and August, respectively.

There were also strong correlations among different climate and remote sensing variables themselves. Strong positive correlation coefficients were observed among the SIF and VIs for both maize and soybean, while LST negatively correlated with other remote sensing variables. VPD also negatively correlated with all the remote sensing variables, except LST with which VPD showed positive correlation coefficients indicating that LST and VPD are good crop stress indicators when crop growth condition is sub-optimal. Compared with VPD, the correlation between Tair and the remote sensing variables were relatively weak.

### 4.2. FFCV of yield prediction performance using only remote sensing variables

We first evaluated the tested yield prediction performance of those models only using remote sensing variables with FFCV out-of-sample validation method. The results for training and testing with data during 2015–2018 and in 2018 only are shown in Fig. 3 and 4, respectively. For maize and soybean yield prediction during 2015–2018, NIRv performed consistently better than other remote sensing variables with the highest $R^2$ and lowest RMSE. The performance of $S\bar{I}F_{OCO2\_005}$ in maize yield prediction was slightly better than NDVI, EVI, and LST, while GOME-2 SIF has the lowest performance with lowest $R^2$ and largest RMSE in crop yield prediction among all the remote sensing variables (Fig. 3). When these models were trained and tested using data in 2018, the performance of using SIF from GOME-2 still showed the lowest performance, while the performances of using other remote sensing variables were quite similar, especially when using nonlinear machine learning algorithms. Overall, we still observed that NIRv, $S\bar{I}F_{OCO2\_005}$, and TROPOMI SIF performed better than other remote sensing
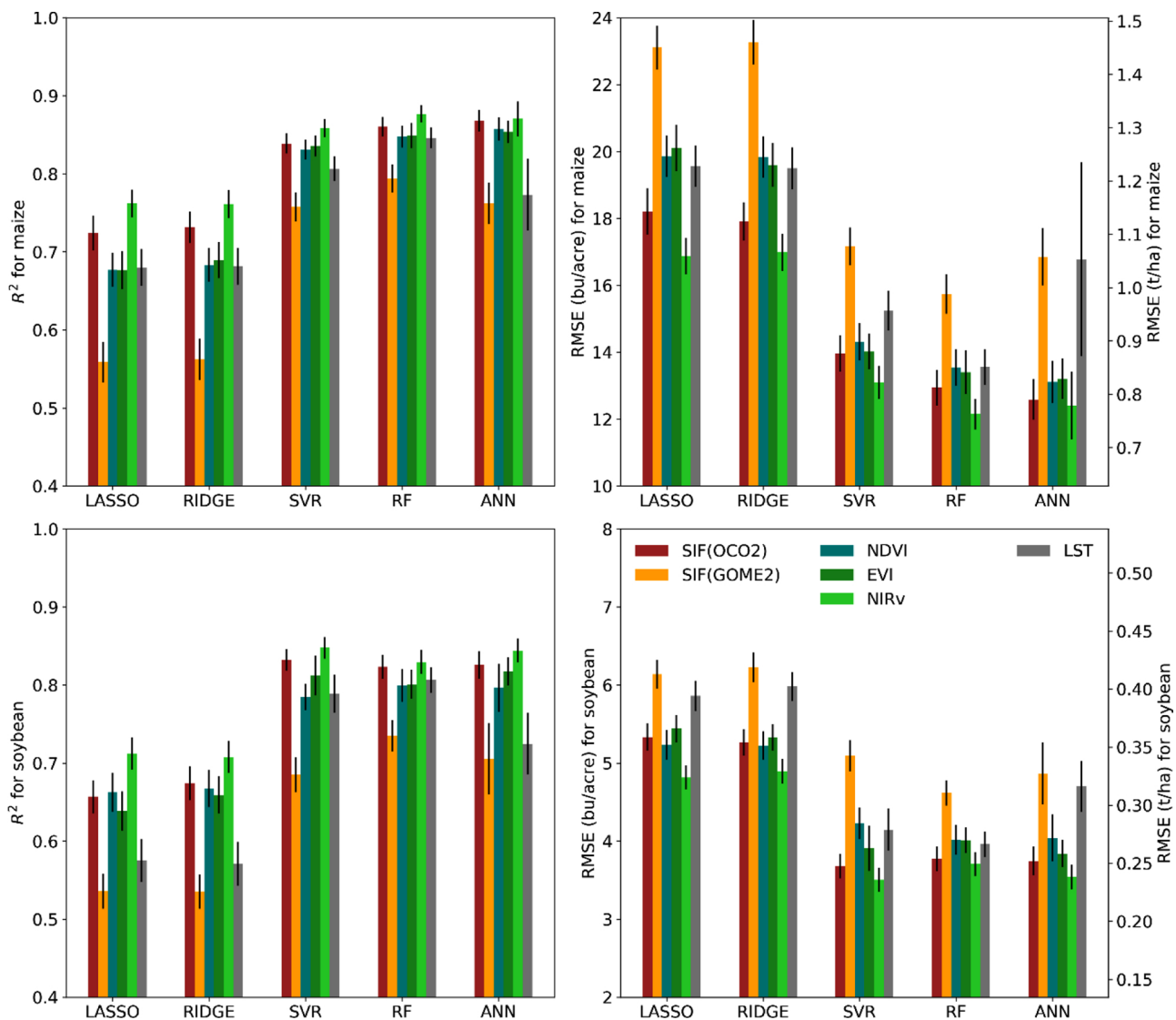
**Fig. 3.** Testing performance of maize (top panels) and soybean (bottom panels) yield prediction using only remote sensing variables and evaluated with five-fold-cross-validation method during 2015–2018. The performance metrics (left panels for $R^2$ and right panels for RMSE) are calculated for 500 random training-testing splits and then both means (filled bars) and standard deviations (error bars) of the metrics are derived. SIF(OCO2) and SIF(GOME2) in the legend represent $\bar{SIF}_{OCO2\_005}$ and GOME-2 SIF products, respectively.

variables (GOME-2 SIF, NDVI, EVI, and LST) in maize and soybean yield prediction (Fig. 4). For maize, NIRv performed consistently better than other remote sensing variables with the highest $R^2$ and lowest RMSE across the five algorithms. For soybean, $\bar{SIF}_{OCO2\_005}$ had a slightly better mean performance compared with NIRv and TROPOMI SIF, but we noted that the performance differences among these three variables were marginal. Results from MAB metric were consistent with the above results from $R^2$ and RMSE (Fig. S1 and S2), i.e., $\bar{SIF}_{OCO2\_005}$, TROPOMI SIF, and NIRv had the lowest MAB among all the remote sensing variables for both maize and soybean yield prediction.

### 4.3. FFCV of yield prediction performance using combined climate and remote sensing variables

We then evaluated the tested yield prediction performance of those models using combined climate and remote sensing variables with FFCV out-of-sample validation method. The results for training and testing with data during 2015–2018 and in 2018 only are shown in Fig. 5 and 6, respectively. For both maize and soybean yield prediction during 2015–2018, NIRv performed best among all the remote sensing

variables with the highest $R^2$ and lowest RMSE (Fig. 5). $\bar{SIF}_{OCO2\_005}$ showed similar performance with EVI or NDVI in yield prediction for both maize and soybean. When the models were trained and tested using data in 2018, NIRv still performed best for maize yield prediction. For soybean yield prediction, NIRv performed best when using linear yield prediction algorithms, while both $\bar{SIF}_{OCO2\_005}$ and VIs had similar performances when using nonlinear yield prediction algorithms. Similar to the results obtained when only using remote sensing variables in crop yield prediction, using SIF from GOME-2 had the lowest performance in crop yield prediction with climate and remote sensing combined models. MAB results showed that $\bar{SIF}_{OCO2\_005}$, TROPOMI SIF, NDVI, EVI, and NIRv had comparable MABs, especially in soybean yield prediction with climate and remote sensing combined models, while GOME-2 SIF and LST had much larger MABs (Fig. S3 and S4).

### 4.4. Forward yield prediction in 2018

The spatial patterns of forward yield prediction in 2018 using different climate and remote sensing combined variables and random forest models are shown in Fig. 7 and 8 for maize and soybean,
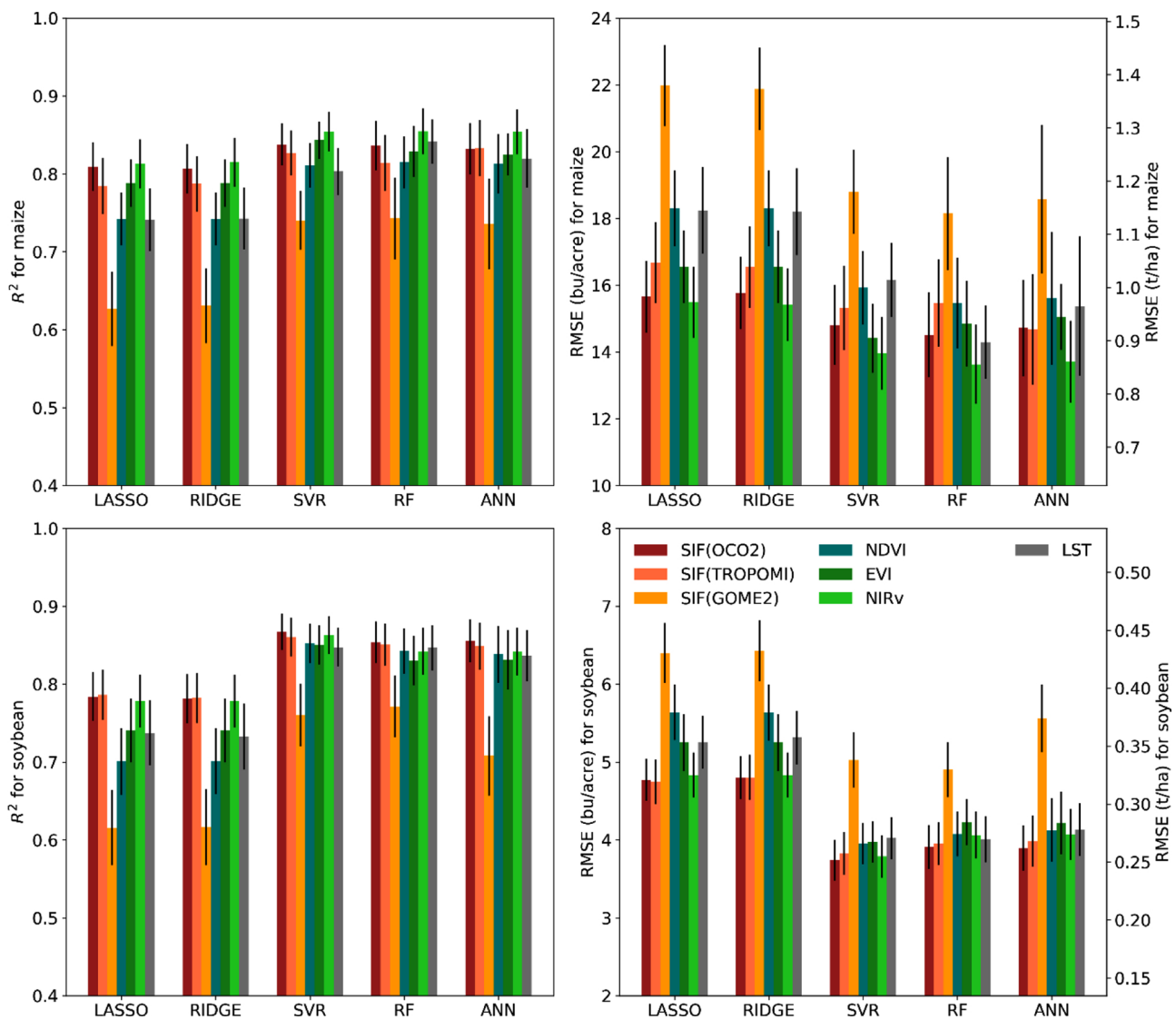
**Fig. 4.** Testing performance of maize (top panels) and soybean (bottom panels) yield prediction using only remote sensing variables and evaluated with five-fold-cross-validation method in 2018. The performance metrics (left panels for $R^2$ and right panels for RMSE) are calculated for 500 random training-testing splits and then both means (filled bars) and standard deviations (error bars) of the metrics are derived. SIF(OCO2), SIF(TROPOMI), and SIF(GOME2) in the legend represent $S\bar{I}F_{OCO2\_005}$, TROPOMI and GOME-2 SIF products, respectively.

respectively. We showed the results from the random forest models because they had the best yield prediction performance as shown in Fig. 3 to Fig. 6. The models were trained using data from 2015-2017. For SIF, we trained the model using $S\bar{I}F_{OCO2\_005}$ data during 2015-2017, while validated the model using both $S\bar{I}F_{OCO2\_005}$ and TROPOMI data in 2018. For both maize and soybean, using $S\bar{I}F_{OCO2\_005}$ and TROPOMI SIF products gave the best yield prediction performances. For example, using $S\bar{I}F_{OCO2\_005}$ in yield prediction in 2018 achieved a $R^2$ of 0.77 and RMSE of 18.11 bu/acre (1.14 t/ha) for maize, and $R^2$ of 0.78 and RMSE of 5.31 bu/acre (0.36 t/ha) for soybean, respectively. Using TROPMI SIF in yield prediction gave similar performance as that using $S\bar{I}F_{OCO2\_005}$ in 2018 with some performance degradation as the models were trained using $S\bar{I}F_{OCO2\_005}$. Using GOME-2 SIF gave the lowest performance in maize yield prediction with an $R^2$ of 0.53 and RMSE of 25.95 bu/acre (1.63 t/ha), while its performance was slightly better than using LST for soybean yield prediction. NIRv performed best among other remote sensing variables besides SIF. Bias distribution in yield prediction showed that using $S\bar{I}F_{OCO2\_005}$, TROPMI SIF, and NIRv could lead to more centralized and narrower bias distributions around

zero compared with using other remote sensing variables for both corn and soybean (Fig. 9).

## 5. Discussion

### 5.1. Potential of using satellite-based SIF products in crop yield prediction

In this study, we demonstrated that using high-resolution SIF products from OCO-2 and TROPOMI could significantly improve the yield prediction performance compared with using coarse-resolution SIF products from GOME-2. This is mainly because higher resolution of $S\bar{I}F_{OCO2\_005}$ and TROPOMI SIF enables better quantification of SIF signals from cropland, and GOME-2 has a lower signal-to-noise ratio. Our study also demonstrated that using high-resolution SIF products from OCO-2 and TROPOMI could bring benefits in crop yield prediction. For example, using SIF products from OCO-2 and TROPOMI achieved the best yield prediction performances for both maize and soybean with either five-fold-cross-validation in 2018 (Fig. 4) or forward prediction in 2018 (Figs. 7, 8, and 9). However, our results also showed that using
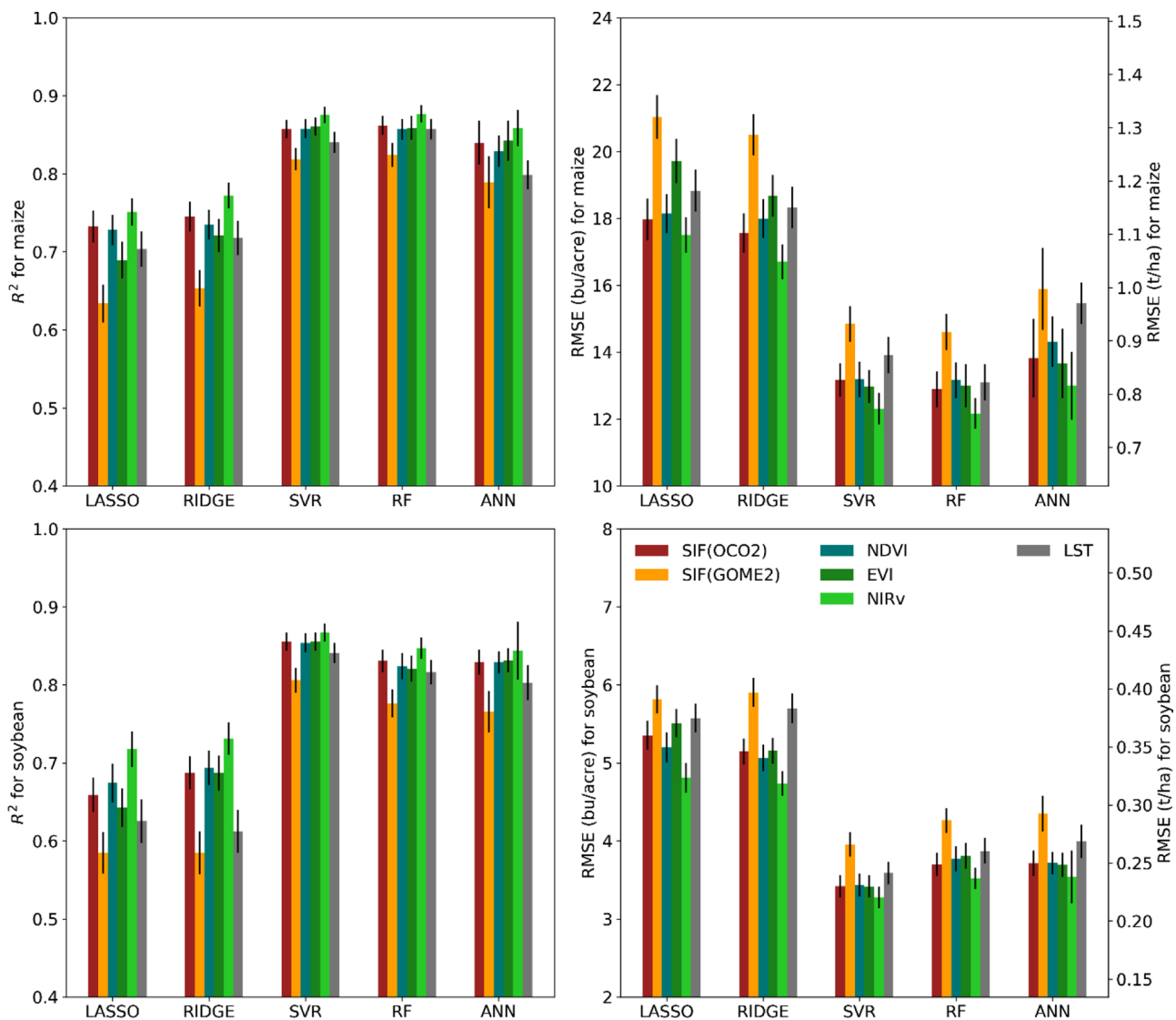
**Fig. 5.** Testing performance of maize (top panels) and soybean (bottom panels) yield prediction using combined climate and remote sensing variables and evaluated with five-fold-cross-validation method during 2015–2018. The performance metrics (left panels for $R^2$ and right panels for RMSE) are calculated for 500 random training-testing splits and then both means (filled bars) and standard deviations (errorbars) of the metrics are derived. SIF(OCO2) and SIF(GOME2) in the legend represents $\overline{SIF}_{OCO2\_005}$ and GOME-2 SIF products, respectively.

current high-resolution SIF products did not guarantee consistently better yield prediction performances than using other remote sensing variables in all the evaluated cases. The relative performances of using different remote sensing variables in yield prediction depended on crop types (maize or soybean), out-of-sample testing methods (five-fold-cross-validation or forward), and length of training data. However, considering that the high-reslution SIF products we used here have a spatial resolution of 5 km while other MODIS-based variables are at finer ($\leq$ 1 km) resolutions, we are still optimistic in the performance of using SIF data for crop yield prediction since higher spatial resolution of SIF data would allow better separation of corn and soybean than the current 5 km data we used in this study.

There are several possible ways that can lead to potential improvement for the yield prediction performance using SIF. Firstly, different ways of using SIF data in building crop yield prediction models may lead to different performances. For example, considering that the SIF signal is integrable over time, we may also use growing season accumulated SIF or maximum SIF in crop yield prediction. Converting satellite-observed SIF (angular SIF from top canopy) into whole-canopy total emitted SIF has been found to better correlated with canopy

photosynthesis (Liu et al. 2019; Yang and van der Tol 2018; Zeng et al. 2019), which may improve the crop yield prediction too. The total emitted SIF from chlorophyll is attenuated by reabsorption and scattering within the leaf and canopy making the observed canopy SIF is a variable fraction of total emitted SIF. The conversion from satellite-observed SIF to total emitted SIF is non-trivial, as we need to estimate the escape ratio, which is determined by sun-sensor geometry, canopy structure, and leaf optical properties. Recent work by Zeng et al. (2019) proposed a practical approach to approximate the escape ratio for near infrared SIF using the NIRv-to-fPAR ratio, making conversion from satellite-observed SIF to total emitted SIF feasible at large scale. To test the performances of these alternative ways of using SIF data in crop yield prediction, we compared the performances of using monthly SIF, growing season maximum and accumulated SIF, and the monthly total emitted SIF from $\overline{SIF}_{OCO2\_005}$ in yield prediction during 2015–2018 (Fig. 10). We used scaling relationships between NDVI and fPAR to calculate fPAR (Gitelson et al. 2014), and subsequently estimated escape ratio and total emitted SIF (Zeng et al. 2019). We observed that using monthly mean SIF actually achieved better performances than other alternative ways of using SIF data for both maize and soybean
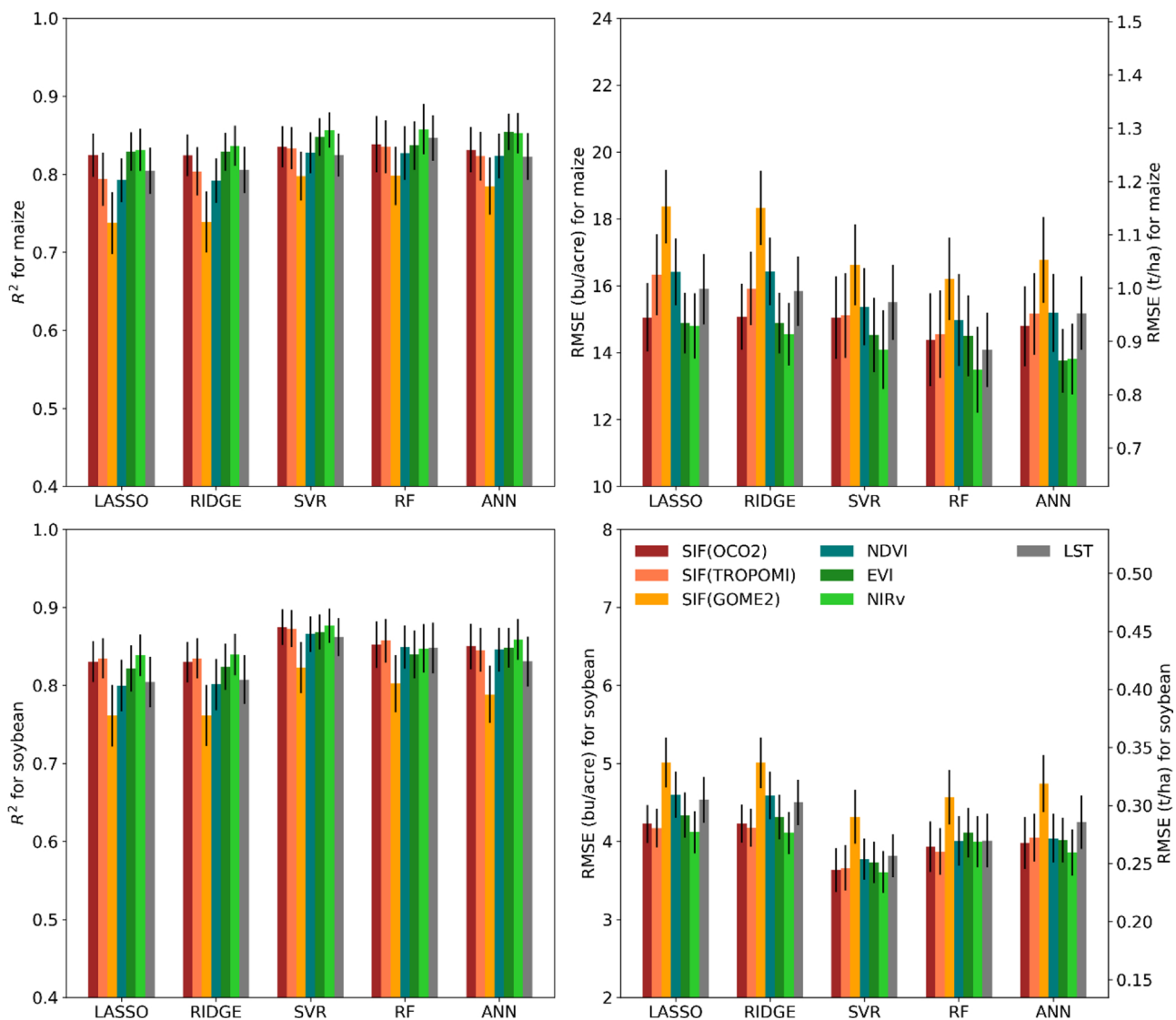
**Fig. 6.** Testing performance of maize (top panels) and soybean (bottom panels) yield prediction using combined climate and remote sensing variables and evaluated with five-fold-cross-validation method in 2018. The performance metrics (left panels for $R^2$ and right panels for RMSE) are calculated for 500 random training-testing splits and then both means (filled bars) and standard deviations (errorbars) of the metrics are derived. SIF(OCO2), SIF(TROPOMI), and SIF(GOME2) in the legend represent $S\bar{I}F_{OCO2\_005}$, TROPOMI and GOME-2 SIF products, respectively.

yield prediction. The reason may be that monthly mean SIF provides more temporal information to capture crop stress at various stages than growing season maximum and accumulated SIF, and there are uncertainties in estimations of both fPAR and escape ratio when deriving the total emitted SIF. Whether using total emitted SIF could lead to improved yield prediction performance when better fPAR estimation and new approximation of escape ratio become available deserves further investigation.

Secondly, we note that the yield prediction performance reported here was only optimized with limited training data. For fair comparison between SIF and other remote sensing predictors, we only used the data during 2015 to 2018 for model training, which may be not enough for operational application aiming at higher performance in yield prediction. With the increase of training data, the yield prediction performance of using SIF products may be further improved. To test this hypothesis, we trained random forest models with different years of training data before 2018 and tested the model performances using data in 2018 (Fig. 11). The models used EVI, NDVI, NIRv and LST as remote sensing variables, all of which are available since early 2000s. These experiments were restricted to six states (Illinois, Indiana, Iowa,

Nebraska, North Dakota, and Wisconsin) with CDL data available since 2003. For maize yield prediction, additional 4 to 6 years (7 to 9 years in total) of training data could further improve the yield prediction performances in 2018. After that, the yield prediction performances were relatively stable, and more training data did not necessarily mean better yield prediction performance any more. For soybean, the change of yield prediction performance with increased years of training data was relatively noisy. These results indicated that the maize yield prediction performances using SIF could be further improved with more training data accumulated. Besides new data from OCO-2/3 and TROPOMI, recent advancements in developing gap-filled OCO-2 SIF products since 2000 (Li and Xiao 2019) and attempts in reconciling inconsistencies among multi-sensor observations in last two decades (Parazoo et al. 2019; Wen et al. 2020) may also help in accumulating more SIF data for training crop yield prediction models, though the uncertainties from the reconstructed SIF data may propagate to the final yield prediction which needs further investigation in future studies.

Thirdly, better quality of future SIF products may further improve the performance in yield prediction. New satellite missions, such as FLuorescence EXplorer (FLEX) (Drusch et al. 2016), can provide SIF
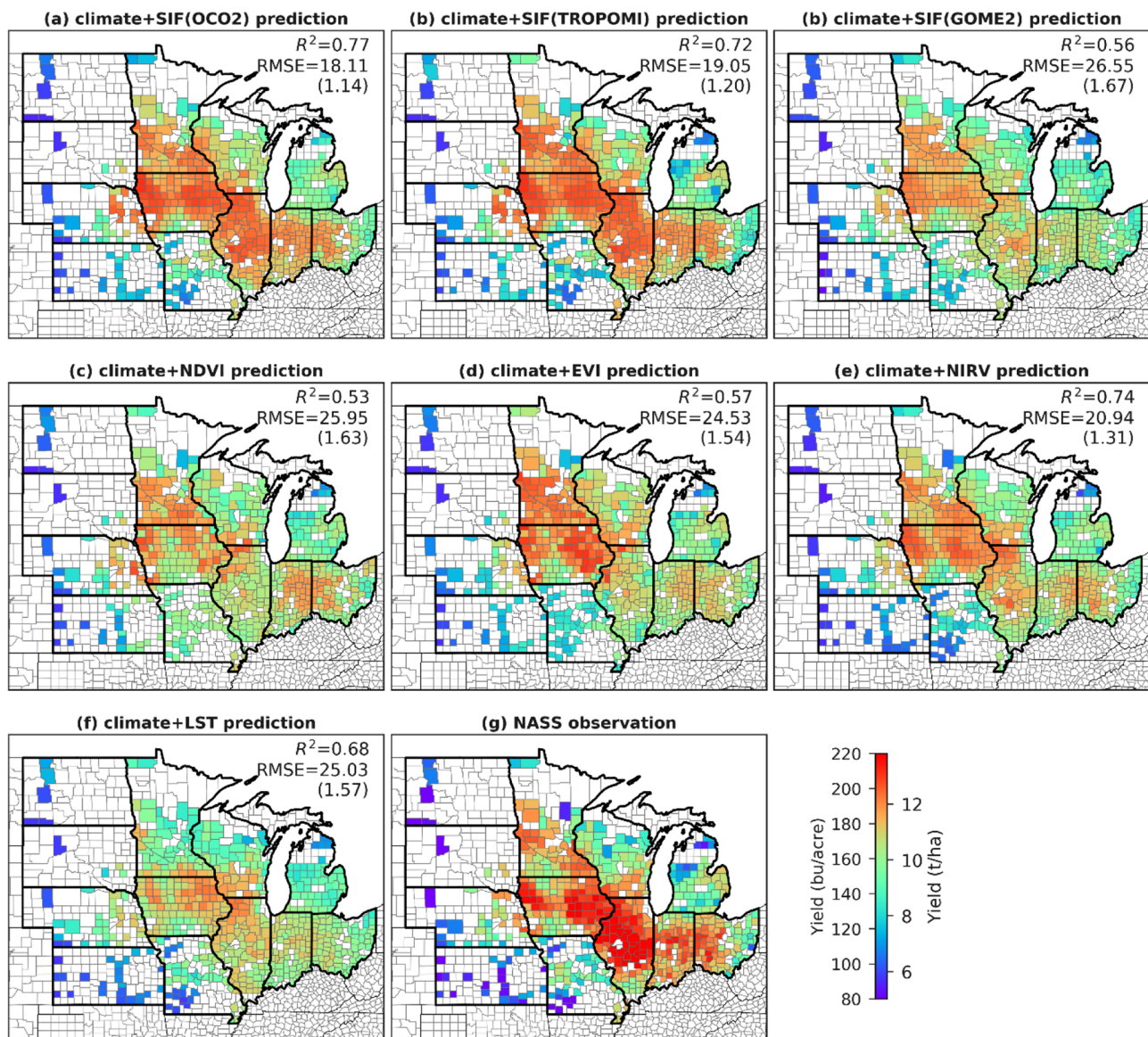
**Fig. 7.** Comparison between the spatial patterns of observed and predicted maize yield in 2018 using random forest model with climate and remote sensing combined variables as inputs. The models were trained using data from 2015-2017. SIF(OCO2), SIF(TROPOMI), and SIF(GOME2) represent $\widetilde{SIF}_{OCO2\_005}$, TROPOMI, and GOME-2 SIF products, respectively. For $\widetilde{SIF}_{OCO2\_005}$ and TROPOMI SIF, we trained the model using $\widetilde{SIF}_{OCO2\_005}$ data during 2015-2017, while validated the model using both $\widetilde{SIF}_{OCO2\_005}$ and TROPOMI data in 2018. RMSE values outside and inside the parentheses are in bu/acre and t/ha, respectively.

products with higher spatial resolutions than existing SIF products. Statistical downscaling also has the potential to further improve the spatial resolution of existing SIF products although previous efforts mainly focused on downscaling the coarse-resolution SIF products, such as those from GOME-2 (Duveiller and Cescatti 2016; Duveiller et al. 2019).

### 5.2. Performance variation among using different remote sensing variables

We also observed differences in the yield prediction performances when using different VIs and LST. Among the tested VIs, NIRv had an overall best performance in predicting maize and soybean yield indicating great potential of using NIRv in crop yield prediction. Compared with traditional remote-sensing-based VIs, NIRv has a more direct physical interpretation as it approximates the proportion of NIR light reflected by vegetation canopy (Badgley et al. 2017). NIRv also minimizes the impacts of soil background and sun-canopy-sensor geometry (Badgley et al. 2019; Badgley et al. 2017). Our study is the first one that used NIRv for crop yield prediction at large scales. NDVI

seemed a good indicator for predicting soybean yield, but not for maize yield. We also found that LST shows better predictability for maize yield than soybean yield (Fig. 11), which may be partly caused by the fact that soybean yield is less sensitive to high temperature and VPD than maize mainly due to relatively higher optimal growth temperature and stable sowing density of soybean over the last two decades (Lobell et al. 2014).

### 5.3. Performance variation among using different machine-learning algorithms

Although we were not aiming to compare the performances of different machine learning algorithms in crop yield estimation in this study, we did see performance difference among the selected five algorithms. Generally, we found the nonlinear algorithms (RF, SVM, and ANN) perform better than the linear algorithms (LASSO and RIDGE). LASSO and RIDGE performed similarly, and the three nonlinear algorithms achieved similar performances in yield prediction for both maize and soybean, which are consistent with previous studies (Cai et al.
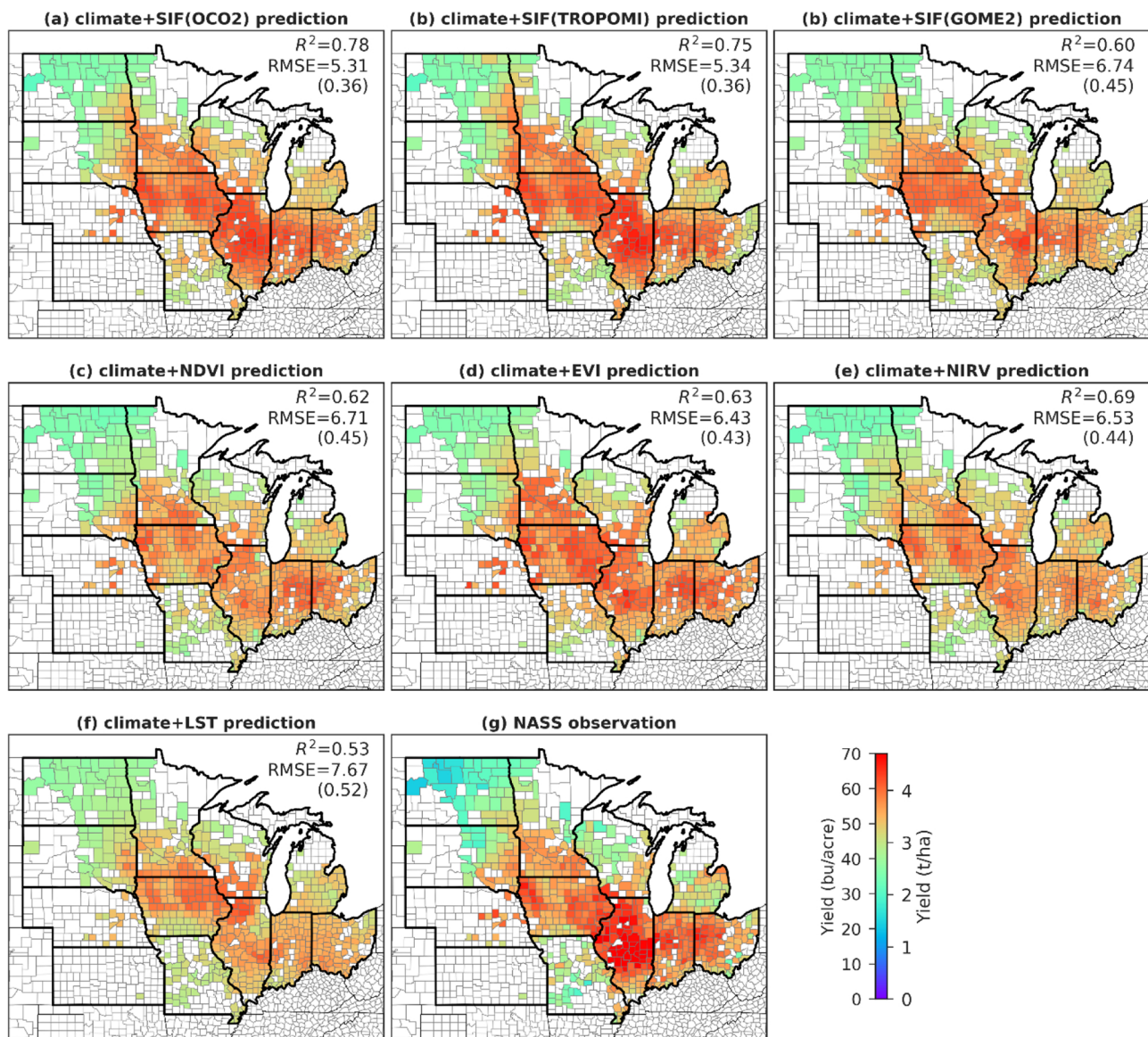
**Fig. 8.** Comparison between the spatial patterns of observed and predicted soybean yield in 2018 using random forest model with climate and remote sensing combined variables as inputs. The models were trained using data from 2015-2017. SIF(OCO2), SIF(TROPOMI), and SIF(GOME2) represent $\bar{SIF}_{OCO2\_005}$, TROPOMI and GOME-2 SIF products, respectively. For $\bar{SIF}_{OCO2\_005}$ and TROPOMI SIF, we trained the model using $\bar{SIF}_{OCO2\_005}$ data during 2015-2017, while validated the model using both $\bar{SIF}_{OCO2\_005}$ and TROPOMI data in 2018. RMSE values outside and inside the parentheses are in bu/acre and t/ha, respectively.

2019). Other advanced machine learning algorithms, such as deep learning algorithms (Oliveira et al. 2018; You et al. 2017), may have the potential to improve the absolute performance in crop yield prediction. However, we expect that the relative performance using SIF and other remote sensing based predictors would not change even when using more advanced algorithms to build crop yield prediction models, which needs further testing.

## 6. Conclusion

With more satellite-based high-resolution SIF products becoming available, there is a need to assess the potential benefits of using these SIF products in operational crop yield prediction. In this study, we evaluated the relative performances of using high-resolution SIF products from OCO-2 and TROPOMI, coarse-resolution SIF product from GOME-2, and MODIS-based VIs (including NDVI, EVI, and NIRv) and LST in predicting maize and soybean yield of the U.S. Midwest. Both remote-sensing-only and climate-remote-sensing-combined yield prediction models were built using five machine-learning algorithms,

including LASSO, RIDGE, SVM, RF and ANN. We found that using high-resolution SIF products from OCO-2 and TROPOMI outperformed using coarse-resolution SIF product from GOME-2 in yield prediction. We also found that using high-resolution SIF products from OCO-2 and TROPOMI gave the best forward predictions for both maize and soybean yields in 2018, indicating great potential of using satellite-based high-resolution SIF products for crop yield prediction. However, using currently available high-resolution SIF products did not guarantee consistently better yield prediction performances than using other satellite-based remote sensing variables in all the evaluated cases, indicating there are still opportunities to improve the quality and resolution of satellite-based high-resolution SIF products. We also found that using NIRv could generally lead to better yield prediction performance than using NDVI, EVI, or LST, and using NIRv could achieve similar or even better yield prediction performance than using the two high-resolution SIF products. These findings indicate that satellite-based high-resolution SIF products could be beneficial in crop yield prediction with more high-resolution and good-quality SIF products accumulated in the future and NIRv is very promising for crop yield
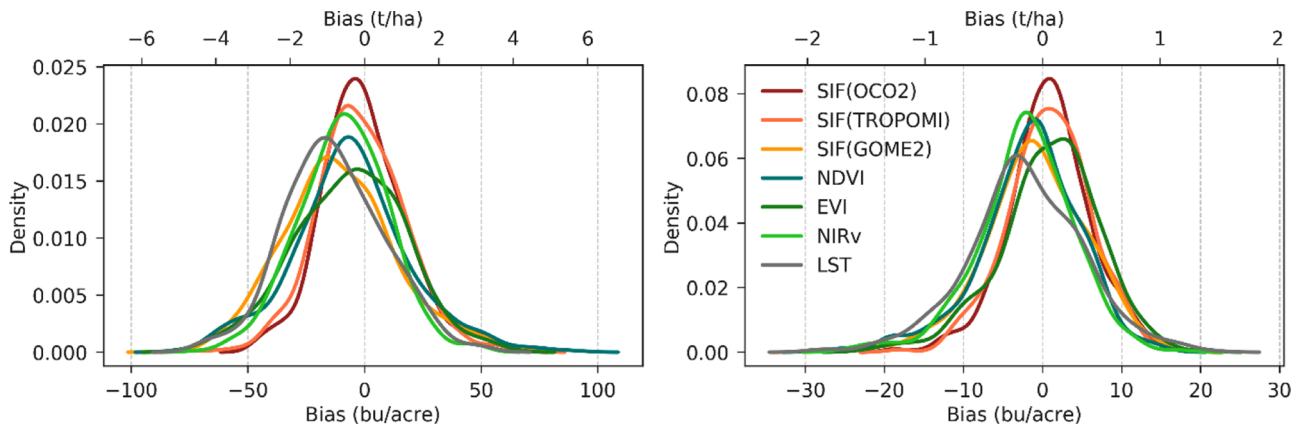
**Fig. 9.** Bias distribution of predicted maize (left) and soybean (right) yield in 2018 using random forest model with climate and remote sensing combined variables as inputs. The models were trained using data from 2015-2017. SIF(OCO2), SIF(TROPOMI), and SIF(GOME2) represent $S\bar{I}F_{OCO2\_005}$, TROPOMI and GOME-2 SIF products, respectively. For $S\bar{I}F_{OCO2\_005}$ and TROPOMI SIF, we trained the model using $S\bar{I}F_{OCO2\_005}$ data during 2015-2017, while validated the model using both $S\bar{I}F_{OCO2\_005}$ and TROPOMI data in 2018.
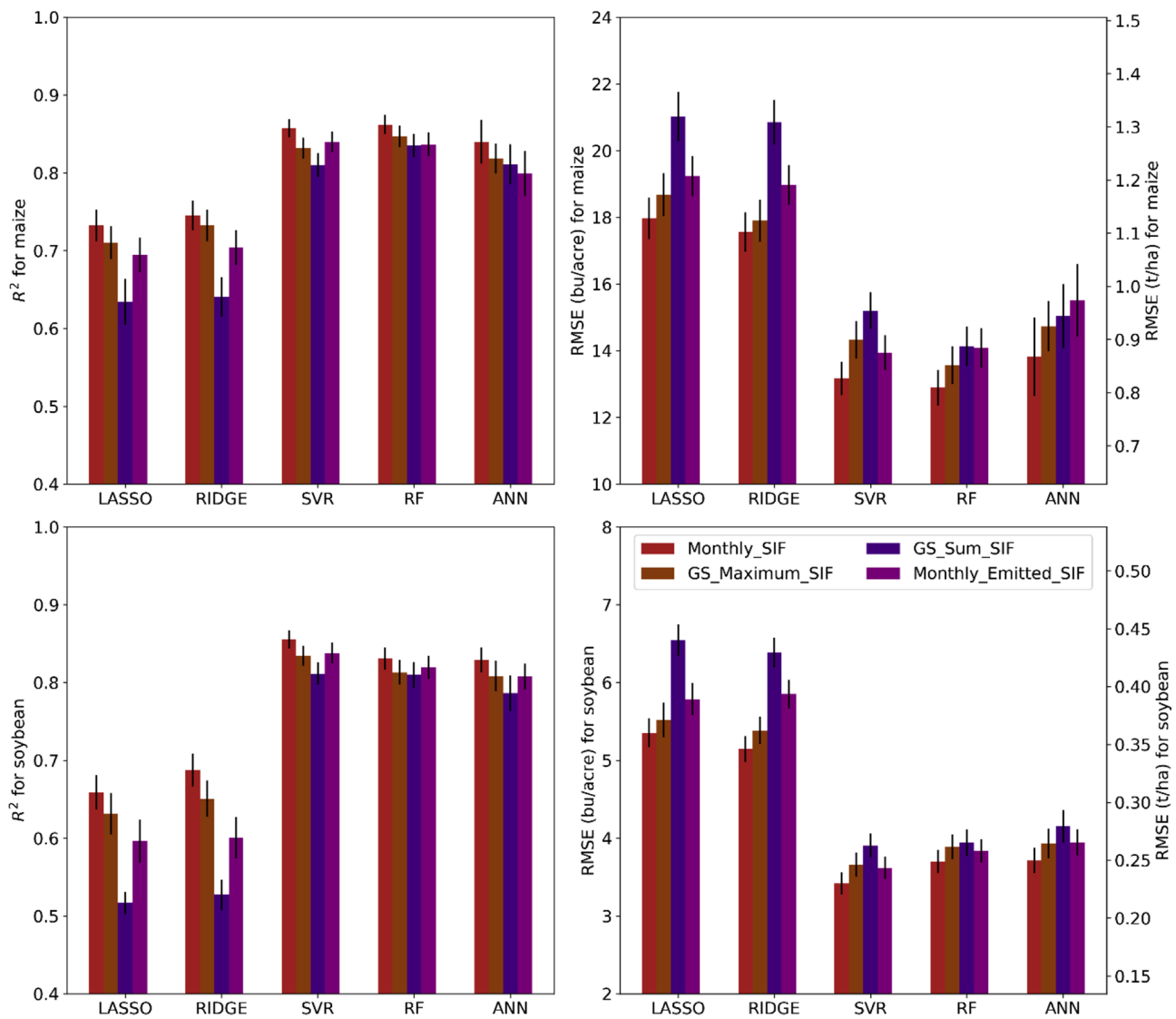


**Fig. 10.** Testing performance of maize (top panels) and soybean (bottom panels) yield prediction using climate variables plus monthly SIF, growing season maximum and accumulated SIF during May to September, or monthly total emitted SIF from $S\bar{I}F_{OCO2\_005}$. The performances were evaluated with five-fold-cross-validation method during 2015–2018. The performance metrics (left panels for $R^2$ and right panels for RMSE) are calculated for 500 random training-testing splits and then both means (filled bars) and standard deviations (errorbars) of the metrics are derived.
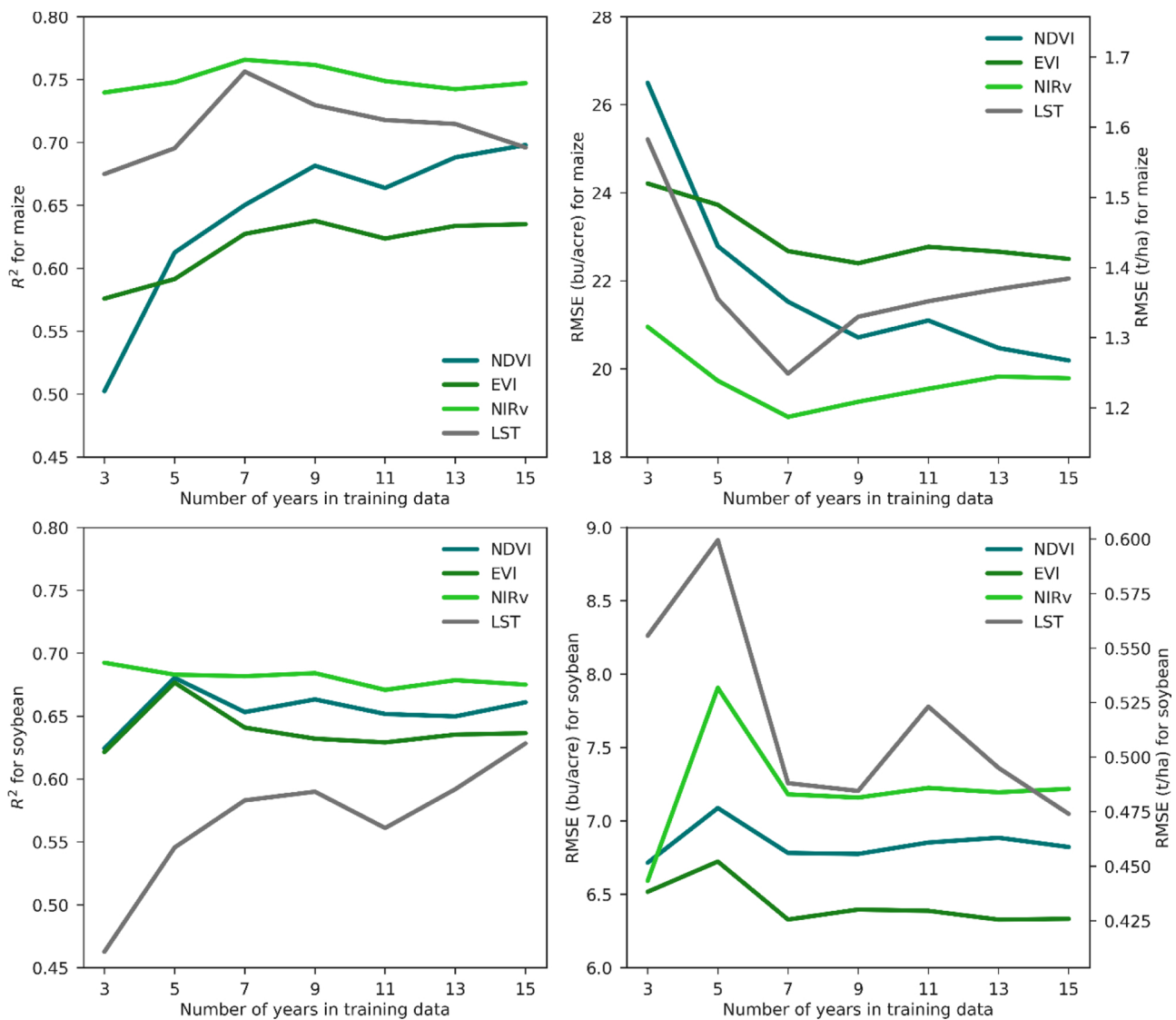
**Fig. 11.** Performance change with increasing years of training data. The random forest models with climate-remote-sensing-combined variables were trained using data before 2018 and validated using data in 2018. The numbers in x-axis represent the number of years before 2018. For example, 3 and 15 in the x-axis mean training data are from 2015 to 2017, and 2003 to 2017, respectively. It has to be noted that model training and validation were conducted over six states with CDL data since 2003, including Illinois, Indiana, Iowa, Nebraska, North Dakota, and Wisconsin.

prediction. To our best knowledge, this study is the first one that compares yield prediction performances of using different SIF products (high-resolution versus coarse-resolution) and using optical-based VIs (including the newly developed NIRv) and thermal-based LST, which can provide insights on developing operational crop yield forecasting system using multi-source remote sensing data. Similar studies outside the U.S. Corn Belt are also needed to assess the performances of using different remote sensing data in crop yield prediction.

**CRediT authorship contribution statement**

**Bin Peng:** Conceptualization, Methodology, Software, Validation, Formal analysis, Investigation, Data curation, Visualization, Writing - original draft, Writing - review & editing. **Kaiyu Guan:** Funding acquisition, Supervision, Project administration, Resources, Conceptualization, Writing - review & editing. **Wang Zhou:** Data curation, Writing - review & editing. **Chongya Jiang:** Writing - review & editing. **Christian Frankenberg:** Data curation, Writing - review & editing. **Ying Sun:** Data curation, Writing - review & editing. **Liyin He:**

Writing - review & editing. **Philipp Köhler:** Data curation, Writing - review & editing.

**Declaration of Competing Interest**

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

**Acknowledgements**

nassgeodata.gmu.edu/CropScape/. MODIS products are available at https://e4ftl01.cr.usgs.gov/. TROPOMI footprint SIF data is available at ftp://fluo.gps.caltech.edu/data/tropomi/ungridded/. $SIF_{OCO2\_005}$ is available at https://cornell.app.box.com/s/cavtg50y80udbdirg022gm5whugmth02. GOME-2 SIF product is available at ftp://fluo.gps.caltech.edu/data/Philipp/GOME-2/. PRISM weather data is available at http://www.prism.oregonstate.edu/.

## Appendix A. Supplementary data

Supplementary material related to this article can be found, in the online version, at doi:https://doi.org/10.1016/j.jag.2020.102126.

## References

Anderson, M.C., Zolin, C.A., Sentelhas, P.C., Hain, C.R., Semmens, K., Tugrul Yilmaz, M., Gao, F., Otkin, J.A., Tetrault, R., 2016. The Evaporative Stress Index as an indicator of agricultural drought in Brazil: An assessment based on crop yield impacts. Remote Sensing of Environment 174, 82–99.

Badgley, G., Anderegg, L.D.L., Berry, J.A., Field, C.B., 2019. Terrestrial Gross Primary Production: Using NIRV to Scale from Site to Globe. Global Change Biology 00, 1–10.

Badgley, G., Field, C.B., Berry, J.A., 2017. Canopy near-infrared reflectance and terrestrial photosynthesis. Science Advances 3, e1602244.

Bastiaanssen, W.G.M., Ali, S., 2003. A new crop yield forecasting model based on satellite measurements applied across the Indus Basin, Pakistan. Agriculture, ecosystems & environment 94, 321–340.

Bolton, D.K., Friedl, M.A., 2013. Forecasting crop yield using remotely sensed vegetation indices and crop phenology metrics. Agricultural and forest meteorology 173, 74–84.

Boryan, C., Yang, Z., Mueller, R., Craig, M., 2011. Monitoring US agriculture: the US department of agriculture, national agricultural statistics service, cropland data layer program. Geocarto International 26, 341–358.

Breiman, L., 2001. Random forests. Mach Learn 45 (5).

Brown, J.N., Hochman, Z., Holzworth, D., Horan, H., 2018. Seasonal climate forecasts provide more definitive and accurate crop yield predictions. Agricultural and forest meteorology 260–261, 247–254.

Cai, Y., Guan, K., Lobell, D., Potgieter, A.B., Wang, S., Peng, J., Xu, T., Asseng, S., Zhang, Y., You, L., Peng, B., 2019. Integrating satellite and climate data to predict wheat yield in Australia using machine learning approaches. Agricultural and forest meteorology 274, 144–159.

Cai, Y., Moore, K., Pellegrini, A., Elhaddad, A., Lessel, J., Townsend, C., Solak, H., Semret, N., 2017. Crop yield predictions - high resolution statistical model for intra-season forecasts applied to corn in the US. Gro Intelligence, Inc.

Chaparro, D., Piles, M., Vall-llossera, M., Camps, A., Konings, A.G., Entekhabi, D., 2018. L-band vegetation optical depth seasonal metrics for crop yield assessment. Remote Sensing of Environment 212, 249–259.

Chipanshi, A., Zhang, Y., Kouadio, L., Newlands, N., Davidson, A., Hill, H., Warren, R., Qian, B., Daneshfar, B., Bedard, F., Reichert, G., 2015. Evaluation of the Integrated Canadian Crop Yield Forecaster (ICCYF) model for in-season prediction of crop yield across the Canadian agricultural landscape. Agricultural and forest meteorology 206, 137–150.

Daly, C., Halbleib, M., Smith, J.I., Gibson, W.P., Doggett, M.K., Taylor, G.H., Curtis, J., Pasteris, P.P., 2008. Physiographically sensitive mapping of climatological temperature and precipitation across the conterminous United States. International Journal of Climatology 28, 2031–2064.

Drusch, M., Moreno, J., Del Bello, U., Franco, R., Goulas, Y., Huth, A., Kraft, S., Middleton, E.M., Miglietta, F., Mohammed, G., 2016. The fluorescence explorer mission concept—ESA's earth explorer 8. Ieee Transactions on Geoscience and Remote Sensing 55, 1273–1284.

Duveiller, G., Cescatti, A., 2016. Spatially downscaling sun-induced chlorophyll fluorescence leads to an improved temporal correlation with gross primary productivity. Remote Sensing of Environment 182, 72–89.

Duveiller, G., Filipponi, F., Walther, S., Köhler, P., Frankenberg, C., Guanter, L., Cescatti, A., 2019. A spatially downscaled sun-induced fluorescence global product for enhanced monitoring of vegetation productivity. Earth Syst. Sci. Data Discuss. 2019, 1–24.

Everingham, Y.L., Muchow, R.C., Stone, R.C., Inman-Bamber, N.G., Singels, A., Bezuidenhout, C.N., 2002. Enhanced risk management and decision-making capability across the sugarcane industry value chain based on seasonal climate forecasts. Agricultural systems 74, 459–477.

Frankenberg, C., Fisher, J.B., Worden, J., Badgley, G., Saatchi, S.S., Lee, J.-E., Toon, G.C., Butz, A., Jung, M., Kuze, A., Yokota, T., 2011. New global observations of the terrestrial carbon cycle from GOSAT: Patterns of plant fluorescence with gross primary productivity. Geophysical Research Letters 38, L17706.

Frankenberg, C., O'Dell, C., Berry, J., Guanter, L., Joiner, J., Köhler, P., Pollock, R., Taylor, T.E., 2014. Prospects for chlorophyll fluorescence remote sensing from the Orbiting Carbon Observatory-2. Remote Sensing of Environment 147, 1–12.

Fu, W.J., 1998. Penalized Regressions: The Bridge versus the Lasso. Journal of computational and graphical statistics 7, 397–416.

Gardner, M.W., Dorling, S., 1998. Artificial neural networks (the multilayer perceptron)—a review of applications in the atmospheric sciences. Atmospheric Environment 32, 2627–2636.

Gitelson, A.A., Peng, Y., Huemmrich, K.F., 2014. Relationship between fraction of radiation absorbed by photosynthesizing maize and soybean canopies and NDVI from remotely sensed data taken at close range and from MODIS 250m resolution data. Remote Sensing of Environment 147, 108–120.

Guan, K., Berry, J., Zhang, Y., Joiner, J., Guanter, L., Badgley, G., Lobell, D.B., 2016. Improving the monitoring of crop productivity using spaceborne solar-induced fluorescence. Global Change Biology 22, 716–726.

Guan, K., Wu, J., Kimball, J.S., Anderson, M.C., Frolking, S., Li, B., Hain, C.R., Lobell, D.B., 2017. The shared and unique values of optical, fluorescence, thermal and microwave satellite data for estimating large-scale crop yields. Remote Sensing of Environment 199, 333–349.

Guanter, L., Zhang, Y., Jung, M., Joiner, J., Voigt, M., Berry, J.A., Frankenberg, C., Huete, A.R., Zarco-Tejada, P., Lee, J.-E., Moran, M.S., Ponce-Campos, G., Beer, C., Camps-Valls, G., Buchmann, N., Gianelle, D., Klumpp, K., Cescatti, A., Baker, J.M., Griffis, T.J., 2014. Global and time-resolved monitoring of crop photosynthesis with chlorophyll fluorescence. Proceedings of the National Academy of Sciences 111, E1327–E1333.

Hansen, J.W., Indeje, M., 2004. Linking dynamic seasonal climate forecasts with crop simulation for maize yield prediction in semi-arid Kenya. Agricultural and forest meteorology 125, 143–157.

He, M., Kimball, S.J., Maneta, P.M., Maxwell, D.B., Moreno, A., Beguería, S., Wu, X., 2018. Regional Crop Gross Primary Productivity and Yield Estimation Using Fused Landsat-MODIS Data. Remote Sensing 10, 372.

Hoerl, A.E., Kennard, R.W., 1970. Ridge regression: Biased estimation for nonorthogonal problems. Technometrics 12, 55–67.

Isengildina-Massa, O., Irwin, S.H., Good, D.L., Gomez, J.K., 2008. The impact of situation and outlook information in corn and soybean futures markets: Evidence from WASDE reports. Journal of Agricultural and Applied Economics 40, 89–103.

Jeong, J.H., Resop, J.P., Mueller, N.D., Fleisher, D.H., Yun, K., Butler, E.E., Timlin, D.J., Shim, K.-M., Gerber, J.S., Reddy, V.R., Kim, S.-H., 2016. Random Forests for Global and Regional Crop Yield Predictions. PLoS ONE 11, e0156571.

Jiang, D., Yang, X., Clinton, N., Wang, N., 2004. An artificial neural network model for estimating crop yields using remotely sensed information. International Journal of Remote Sensing 25, 1723–1732.

Johnson, D.M., 2014. An assessment of pre- and within-season remotely sensed variables for forecasting corn and soybean yields in the United States. Remote Sensing of Environment 141, 116–128.

Joiner, J., Guanter, L., Lindstrot, R., Voigt, M., Vasilkov, A., Middleton, E., Huemmrich, K., Yoshida, Y., Frankenberg, C., 2013. Global monitoring of terrestrial chlorophyll fluorescence from moderate-spectral-resolution near-infrared satellite measurements: methodology, simulations, and application to GOME-2. Atmospheric Measurement Techniques 6, 2803–2823.

Jones, J.W., Antle, J.M., Basso, B., Boote, K.J., Conant, R.T., Foster, I., Godfray, H.C.J., Herrero, M., Howitt, R.E., Janssen, S., Keating, B.A., Munoz-Carpena, R., Porter, C.H., Rosenzweig, C., Wheeler, T.R., 2017. Toward a new generation of agricultural system data, models, and knowledge products: State of agricultural systems science. Agricultural systems 155, 269–288.

Jones, J.W., Hoogenboom, G., Porter, C.H., Boote, K.J., Batchelor, W.D., Hunt, L.A., Wilkens, P.W., Singh, U., Gijsman, A.J., Ritchie, J.T., 2003. The DSSAT cropping system model. European Journal of Agronomy 18, 235–265.

Köhler, P., Frankenberg, C., Magney, T.S., Guanter, L., Joiner, J., Landgraf, J., 2018. Global retrievals of solar induced chlorophyll fluorescence with TROPOMI: first results and inter-sensor comparison to OCO-2. Geophysical Research Letters 45 (10), 456–463.

Köhler, P., Guanter, L., Joiner, J., 2015. A linear method for the retrieval of sun-induced chlorophyll fluorescence from GOME-2 and SCIAMACHY data. Atmos. Meas. Tech. 8, 2589–2608.

Legler, D.M., Bryant, K.J., O'Brien, J.J., 1999. Impact of ENSO-Related Climate Anomalies on Crop Yields in the U.S. Climatic Change 42, 351–375.

Li, X., Xiao, J., 2019. A Global, 0.05-Degree Product of Solar-Induced Chlorophyll Fluorescence Derived from OCO-2, MODIS, and Reanalysis Data. Remote Sensing 11, 517.

Li, Y., Guan, K., Yu, A., Peng, B., Zhao, L., Li, B., Peng, J., 2019. Toward building a transparent statistical model for improving crop yield prediction: Modeling rainfed corn in the U.S. Field Crops Research 234, 55–65.

Liu, X., Guanter, L., Liu, L., Damm, A., Malenovský, Z., Rascher, U., Peng, D., Du, S., Gastellu-Etchegorry, J.-P., 2019. Downscaling of solar-induced chlorophyll fluorescence from canopy level to photosystem level using a random forest model. Remote Sensing of Environment 203, 110772.

Lobell, D.B., Burke, M.B., 2010. On the use of statistical models to predict crop yield responses to climate change. Agricultural and forest meteorology 150, 1443–1452.

Lobell, D.B., Roberts, M.J., Schlenker, W., Braun, N., Little, B.B., Rejesus, R.M., Hammer, G.L., 2014. Greater Sensitivity to Drought Accompanies Maize Yield Increase in the U.S. Midwest. Science 344, 516–519.

Lobell, D.B., Thau, D., Seifert, C., Engle, E., Little, B., 2015. A scalable satellite-based crop yield mapper. Remote Sensing of Environment 164, 324–333.

Magney, T.S., Frankenberg, C., Fisher, J.B., Sun, Y., North, G.B., Davis, T.S., Kornfeld, A., Siebke, K., 2017. Connecting active to passive fluorescence with photosynthesis: a method for evaluating remote sensing measurements of Chl fluorescence. New Phytologist 215, 1594–1608.

Newlands, N.K., Zamar, D.S., Kouadio, L.A., Zhang, Y., Chipanshi, A., Potgieter, A., Toure, S., Hill, H.S.J., 2014. An integrated, probabilistic model for improved seasonal forecasting of agricultural crop yield under environmental uncertainty. Frontiers in Environmental Science 2, 1–21.

Oliveira, I., Cunha, R.L., Silva, B., Netto, M.A., 2018. A Scalable Machine Learning System for Pre-Season Agriculture Yield Forecast. arXiv preprint arXiv 1806, 09244.

Parazoo, N.C., Bowman, K., Frankenberg, C., Lee, J.-E., Fisher, J.B., Worden, J., Jones, D.B.A., Berry, J., Collatz, G.J., Baker, I.T., Jung, M., Liu, J., Osterman, G., O'Dell, C., Sparks, A., Butz, A., Guerlet, S., Yoshida, Y., Chen, H., Gerbig, C., 2013. Interpreting seasonal changes in the carbon balance of southern Amazonia using measurements of XCO2 and chlorophyll fluorescence from GOSAT. Geophysical Research Letters 40, 2829–2833.

Parazoo, N.C., Frankenberg, C., Köhler, P., Joiner, J., Yoshida, Y., Magney, T., Sun, Y., Yadav, V., 2019. Towards a harmonized long-term spaceborne record of far-red solar induced fluorescence. Journal of Geophysical Research: Biogeosciences.

Peng, B., Guan, K., Chen, M., Lawrence, D.M., Pokhrel, Y., Suyker, A., Arkebauer, T., Lu, Y., 2018a. Improving maize growth processes in the community land model: Implementation and evaluation. Agricultural and forest meteorology 250–251, 64–89.

Peng, B., Guan, K., Pan, M., Li, Y., 2018b. Benefits of seasonal climate prediction and satellite data for forecasting US maize yield. Geophysical Research Letters 45, 9662–9671.

Peng, B., Guan, K., Tang, J., Ainsworth, E.A., Asseng, S., Bernacchi, C.J., Cooper, M., Delucia, E.H., Elliott, J.W., Ewert, F., Grant, R.F., Gustafson, D.I., Hammer, G.L., Jin, Z., Jones, J.W., Kimm, H., Lawrence, D.M., Li, Y., Lombardozzi, D.L., Marshall-Colon, A., Messina, C.D., Ort, D.R., Schnable, J.C., Vallejos, C.E., Wu, A., Yin, X., Zhou, W., 2020. Towards a multiscale crop modelling framework for climate change adaptation assessment. Nature Plants 6, 338–348.

Phillips, J., Rajagopalan, B., Cane, M., Rosenzweig, C., 1999. The role of ENSO in determining climate and maize yield variability in the US cornbelt. International Journal of Climatology 19, 877–888.

Potgieter, A.B., Hammer, G.L., Butler, D., 2002. Spatial and temporal patterns in Australian wheat yield and their relationship with ENSO. Australian Journal of Agricultural Research 53, 77–89.

Qian, B., De Jong, R., Warren, R., Chipanshi, A., Hill, H., 2009. Statistical spring wheat yield forecasting for the Canadian prairie provinces. Agricultural and forest meteorology 149, 1022–1031.

Rosenzweig, C., Jones, J.W., Hatfield, J.L., Ruane, A.C., Boote, K.J., Thorburn, P., Antle, J.M., Nelson, G.C., Porter, C., Janssen, S., Asseng, S., Basso, B., Ewert, F., Wallach, D., Baigorria, G., Winter, J.M., 2013. The Agricultural Model Intercomparison and Improvement Project (AgMIP): Protocols and pilot studies. Agricultural and forest meteorology 170, 166–182.

Shelia, V., Hansen, J., Sharda, V., Porter, C., Aggarwal, P., Wilkerson, C.J., Hoogenboom, G., 2019. A Multi-scale and Multi-model Gridded Framework for Forecasting Crop Production, Risk Analysis, and Climate Change Impact Studies. Environmental Modelling & Software 115, 144–154.

Shiga, Y.P., Tadić, J.M., Qiu, X., Yadav, V., Andrews, A.E., Berry, J.A., Michalak, A.M., 2018. Atmospheric CO2 observations reveal strong correlation between regional net biospheric carbon uptake and solar induced chlorophyll fluorescence. Geophysical Research Letters 45, 1122–1132.

Smola, A.J., Schölkopf, B., 2004. A tutorial on support vector regression. Statistics and Computing 14, 199–222.

Specht, D.F., 1991. A general regression neural network. IEEE transactions on neural networks 2, 568–576.

Sun, Y., Frankenberg, C., Jung, M., Joiner, J., Guanter, L., Köhler, P., Magney, T., 2018. Overview of Solar-Induced chlorophyll Fluorescence (SIF) from the Orbiting Carbon Observatory-2: Retrieval, cross-mission comparison, and global monitoring for GPP. Remote Sensing of Environment 209, 808–823.

Suykens, J.A., Vandewalle, J., 1999. Least squares support vector machine classifiers. Neural processing letters 9, 293–300.

Tibshirani, R., 1996. Regression shrinkage and selection via the lasso. Journal of the Royal Statistical Society. Series B (Methodological) 267–288.

Wen, J., Köhler, P., Duveiller, G., Parazoo, N.C., Magney, T.S., Hooker, G., Yu, L., Chang, C.Y., Sun, Y., 2020. A framework for harmonizing multiple satellite instruments to generate a long-term global high spatial-resolution solar-induced chlorophyll fluorescence (SIF). Remote Sensing of Environment 239, 111644.

Yang, P., van der Tol, C., 2018. Linking canopy scattering of far-red sun-induced chlorophyll fluorescence with reflectance. Remote Sensing of Environment 209, 456–467.

Yang, Y., Anderson, M.C., Gao, F., Wardlow, B., Hain, C.R., Otkin, J.A., Alfieri, J., Yang, Y., Sun, L., Dulaney, W., 2018. Field-scale mapping of evaporative stress indicators of crop yield: An application over Mead, NE, USA. Remote Sensing of Environment 210, 387–402.

You, J., Li, X., Low, M., Lobell, D., Ermon, S., 2017. Deep Gaussian Process for Crop Yield Prediction Based on Remote Sensing Data. AAAI, pp. 4559–4566.

Yu, L., Wen, J., Chang, C.Y., Frankenberg, C., Sun, Y., 2018. High Resolution Global Contiguous Solar-Induced Chlorophyll Fluorescence (SIF) of Orbiting Carbon Observatory-2 (OCO-2). Geophysical Research Letters 46, 1449–1458.

Zeng, Y., Badgley, G., Dechant, B., Ryu, Y., Chen, M., Berry, J.A., 2019. A practical approach for estimating the escape ratio of near-infrared solar-induced chlorophyll fluorescence. Remote Sensing of Environment, 111209.

Zhang, Y., Joiner, J., Alemohammad, S.H., Zhou, S., Gentine, P., 2018. A global spatially contiguous solar-induced fluorescence (CSIF) dataset using neural networks. Biogeosciences 15, 5779–5800.