

# HIGH SPATIAL RESOLUTION FOR CROP YIELD PREDICTION IN LARGE FARMING SYSTEMS: A NECESSITY OR ADDITIONAL OVERHEAD

*Stella Ofori-Ampofo<sup>1</sup>, Ridvan Salih Kuzu<sup>2</sup>, Xiao Xiang Zhu<sup>1</sup>*

<sup>1</sup>Technical University of Munich

<sup>2</sup>German Aerospace Center (DLR)

## ABSTRACT

Concerning county-level yield predictions in large farming systems, relying on coarse-resolution satellite images is customary. However, these images lack sufficient textural detail to accurately summarise spatial information. Our work evaluates the advantage of enhanced spatial resolution by conducting a comparative analysis between coarse-resolution, high-temporal-frequency MODIS data and relatively high-resolution, low-temporal-frequency Landsat data for predicting corn yield in the USA. We benchmark this comparison using several models in a spatial versus non-spatial input data case. According to our findings, incorporating high-spatial resolution in this context did not yield any significant benefits, as it comes at the cost of reduced temporal revisit. Root mean square error of about 13 bushels per acre can be achieved 2-4 weeks before harvest for less-intense drought years. Extreme drought-struck years are, however difficult to anticipate.

*Index Terms*— crop yield prediction, machine learning, convolutional neural network, recurrent neural network.

## 1. INTRODUCTION

Achieving zero hunger remains challenging due to climate extremes, economic shocks, and conflicts [1]. Therefore, precise and timely prediction of crop yields is crucial for decisions to control sudden food shortages. In this regard, open-access satellite data offers a wealth of spatio-temporal information that can be effectively utilized by integration of remote sensing and artificial intelligence. This combination has demonstrated great potential in various agricultural applications, such as crop condition assessment [2], vulnerability evaluation, and yield estimation [3, 4, 5].

There are studies utilizing different combinations of variables, such as satellite surface reflectance, vegetation indices, climate data [6, 7], as well as soil properties and farm management information [3], employing various machine learning (ML) and deep learning (DL) approaches. In large-scale farming systems, such as those found in the United States, coarse-resolution satellite images, often acquired from the Moderate Resolution Imaging Spectroradiometer (MODIS),

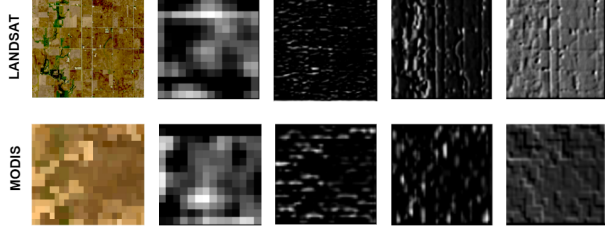
are commonly utilized. To extract spatio-spectral and temporal features, volumetric convolutional neural networks (CNN) are typically employed [4, 7]. Alternatively, two-dimensional CNNs combined with recurrent neural networks (RNN) are used to model the temporal component, which generally outperforms individual CNN or RNN models [6, 8]. However, coarse-resolution images may lack sufficient textural information to capture representative spatio-spectral features. An advantage of higher spatial resolution is its ability to capture fine-grained variations, enabling the analysis of spatial patterns that can provide insights into crop textural variations. Their added benefit has been explored in similar tasks but for smaller farming systems, with shorter periods using synthetic data and without leveraging spatial insights [5, 9].

In this study, we present findings on county-level corn yield prediction by comparing coarse spatial yet high temporal resolution imagery (MODIS-MOD09A1.061) to higher spatial yet lower temporal resolution imagery (Landsat-8). The spatial resolution of MODIS and Landsat-8 is 500 and 30 meters with a temporal frequency of 16 and 8-day respectively. These surface reflectance products are complemented with weather variables and our baseline approach involves encoding without spatial information (pixel averages), and spatial order (time series as histograms) [4].

## 2. DATA

Corn production dominates in the mid-western US; hence, we limit our analysis to the top five production states: Iowa, Illinois, Indiana, Nebraska, and Minnesota. On average, farm sizes in these selected states are over 250 acres [10], which encourages using of coarse-resolution data for yield applications. Google Earth Engine [11] is used to collect surface reflectances and Daymet climate variables (precipitation, minimum, and maximum temperature) for the usual planting and harvesting period (April to September).

MODIS contains seven spectral reflectance bands within the visible, near infra-red and short-wave-infra-red regions. For Landsat, we retain six out of eleven bands to align with the spectral range of MODIS. The daily granularity of the Daymet variables are reduced to coincide with the 8-day and 16-day observation window of MODIS and Landsat. We fur-



**Fig. 1.** Feature maps derived from convolution operations on Landsat and MODIS

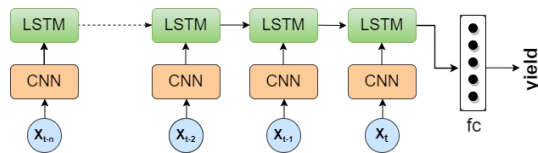
then compute normalized difference water (NDWI) and vegetation indices (NDVI) to introduce expert-knowledge features and mask non-corn pixels using crop data layer from the United States Department of Agriculture (USDA).

Two main data input formats are investigated: (i) a satellite time series  $X \in R^{T \times C}$  created by averaging pixels within each county, (ii) a patch-level input of shape  $X \in R^{T \times C \times H \times W}$ . For patch inputs, we subsequently transform satellite image time series into compact three-dimensional histograms of binarized pixel counts  $X \in R^{b \times T \times C}$  [4] where  $b$  is the number of bins and  $T$ ,  $C$ ,  $H$ , and  $W$  are the time series length, number of channels, image height and width respectively. The target variable, county-level average yield (in bushel per acre) is obtained from USDA, National Agricultural Statistics Service from period 2010 to 2021.

### 3. EXPERIMENTAL SETTING

We design a joint CNN-LSTM architecture (Figure 2) to simultaneously learn spatio-spectral and temporal patterns for patch inputs. The architecture consist of three blocks of time-distributed convolution layers, each followed by a batch normalization, rectified linear unit (ReLU) activation and max-pooling. The extracted spatio-spectral features are passed to an LSTM to learn temporal dependency. To avoid feeding the whole patch of a county (computationally challenging for Landsat data), five sub-patches are randomly sampled per county as inputs.

In the absence of spatial information, Random Forest (RF) is applied. Likewise, the CNN blocks shown in Figure 2 are removed to focus solely on modeling temporal sequences using LSTM. In the case of histogram inputs, two-dimensional convolutions are employed, following a similar approach as described in [4].



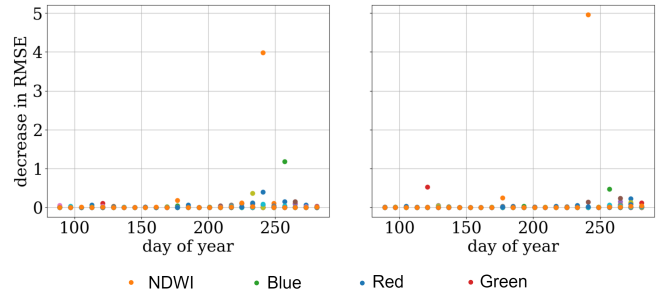
**Fig. 2.** Spatio-spectral and temporal feature extraction using a CNN-LSTM

Our experiments consist of two scenarios. *Scenario-1* considers only MODIS data from 2010 to 2021 across five states, increasing the overall training sample size. Scenario-1 includes observations from diverse agro-climatic zones and enables the comparison of different input data formats for test year 2020 and drought year 2012. [12]. Due to the absence of Landsat-8 data for 2012, *Scenario-2* strictly compares MODIS and Landsat data for Iowa state only. It considers observations from 2013 to 2021, with 2020 serving as an independent test set. The validation set comprises 25% of the training samples in all setups.

## 4. RESULTS AND DISCUSSION

### 4.1. Feature importance across time using Scenario-1

Using permutation-based importance on RF model (Figure 3), we observe for both cases of drought and non-drought the significance of NDWI during senescence. RMSE is reduced by an average of 4-5 bu/acre without the NDWI feature. The effect of the blue band, indicative of water presence, is seen to have a double impact during drought years. However, constructing a RF model using only the top significant features, NDWI and blue features, increases the RMSE to 44 and 25 bu/acre for 2012 and 2020.



**Fig. 3.** Scenario-1: Temporal feature relevance using permutation importance on test sets. Drought-struck year 2012 (left), non-drought year 2020 (right). Legend is limited to the most significant features

### 4.2. Comparing performance for drought-struck and non-drought years

Drought year (2012) was more challenging to predict as it exhibits above-normal weather conditions during crop development, and presents an unprecedented sample in our time series. From Table 1, all models are competitive without extreme drought impacts but the LSTM explained over 50% of the variance for both test years. With extensive training samples (introducing geographical variability), we observe a general improvement in performance for the year 2020 compared to the same year in Scenario-2, where only Iowa state is considered. The latter generalises poorly to an unseen year. The relatively poor performance of CNN-LSTM in Scenario-1 can

**Table 1.** Performance metrics use only MODIS data (Scenario-1) over the top-five corn-growing states. Training data consists of satellite time series from 2010 to 2021 except for 2012 and 2020, which are independent test sets

| Model             | Validation                 |                             |                             | Test (2012)                  |                              |                             | Test (2020)                |                              |                             |
|-------------------|----------------------------|-----------------------------|-----------------------------|------------------------------|------------------------------|-----------------------------|----------------------------|------------------------------|-----------------------------|
|                   | MAPE                       | RMSE                        | R <sup>2</sup>              | MAPE                         | RMSE                         | R <sup>2</sup>              | MAPE                       | RMSE                         | R <sup>2</sup>              |
| RF (M)            | 4.63 <sup>±0.04</sup>      | 10.07 <sup>±0.06</sup>      | 0.86 <sup>±0.01</sup>       | 45.86 <sup>±1.1</sup>        | 39.90 <sup>±0.91</sup>       | 0.03 <sup>±0.05</sup>       | 5.97 <sup>±0.15</sup>      | 13.49 <sup>±0.27</sup>       | 0.62 <sup>±0.02</sup>       |
| Histogram-CNN (M) | 4.50 <sup>±0.0</sup>       | 9.72 <sup>±0.23</sup>       | 0.88 <sup>±0.01</sup>       | 34.40 <sup>±0.01</sup>       | 31.85 <sup>±0.68</sup>       | 0.40 <sup>±0.03</sup>       | 6.12 <sup>±0.0</sup>       | 13.66 <sup>±0.25</sup>       | 0.62 <sup>±0.01</sup>       |
| LSTM (M)          | <b>4.48<sup>±0.0</sup></b> | <b>9.86<sup>±0.26</sup></b> | <b>0.87<sup>±0.01</sup></b> | <b>30.70<sup>±0.01</sup></b> | <b>27.75<sup>±0.79</sup></b> | <b>0.53<sup>±0.03</sup></b> | <b>5.71<sup>±0.0</sup></b> | <b>12.78<sup>±0.80</sup></b> | <b>0.66<sup>±0.04</sup></b> |
| CNN-LSTM (M)      | 5.78 <sup>±0.01</sup>      | 12.84 <sup>±2.08</sup>      | 0.79 <sup>±0.06</sup>       | 55.08 <sup>±0.01</sup>       | 50.31 <sup>±0.87</sup>       | -0.49 <sup>±0.05</sup>      | 7.32 <sup>±0.0</sup>       | 16.30 <sup>±0.08</sup>       | 0.46 <sup>±0.01</sup>       |

**Table 2.** Comparison of Landsat (L) and MODIS (M) in Iowa state (Scenario-2). Training data consists of satellite time series from 2013 to 2021 except for 2020 which is independently used as a test set

| Model             | Validation                  |                             |                             | Test (2020)                 |                              |                             |
|-------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|------------------------------|-----------------------------|
|                   | MAPE                        | RMSE                        | R <sup>2</sup>              | MAPE                        | RMSE                         | R <sup>2</sup>              |
| RF (L)            | 4.07 <sup>±0.17</sup>       | 10.06 <sup>±0.38</sup>      | 0.76 <sup>±0.02</sup>       | 10.34 <sup>±0.30</sup>      | 21.61 <sup>±0.47</sup>       | -0.82 <sup>±0.08</sup>      |
| RF (M)            | 3.45 <sup>±0.23</sup>       | 8.04 <sup>±0.52</sup>       | 0.86 <sup>±0.02</sup>       | 11.98 <sup>±0.87</sup>      | 24.24 <sup>±1.30</sup>       | -1.18 <sup>±0.23</sup>      |
| Histogram-CNN (L) | 4.46 <sup>±0.0</sup>        | 10.68 <sup>±0.05</sup>      | 0.73 <sup>±0.01</sup>       | 8.43 <sup>±0.01</sup>       | 18.04 <sup>±0.85</sup>       | -0.30 <sup>±0.12</sup>      |
| Histogram-CNN (M) | 4.19 <sup>±0.0</sup>        | 9.95 <sup>±0.81</sup>       | 0.78 <sup>±0.02</sup>       | 7.73 <sup>±0.01</sup>       | 16.66 <sup>±1.90</sup>       | -0.03 <sup>±0.23</sup>      |
| LSTM (L)          | 4.25 <sup>±0.0</sup>        | 10.02 <sup>±0.2</sup>       | 0.79 <sup>±0.01</sup>       | 10.69 <sup>±0.01</sup>      | 22.27 <sup>±1.42</sup>       | -0.94 <sup>±0.24</sup>      |
| LSTM (M)          | <b>3.30<sup>±0.00</sup></b> | <b>7.71<sup>±0.53</sup></b> | <b>0.87<sup>±0.01</sup></b> | <b>6.98<sup>±0.01</sup></b> | <b>15.36<sup>±1.75</sup></b> | <b>0.12<sup>±0.21</sup></b> |
| CNN-LSTM (L)      | 9.50 <sup>±0.0</sup>        | 21.14 <sup>±0.86</sup>      | 0.01 <sup>±0.02</sup>       | 10.4 <sup>±0.01</sup>       | 21.53 <sup>±2.59</sup>       | -0.71 <sup>±0.41</sup>      |
| CNN-LSTM (M)      | 8.01 <sup>±0.0</sup>        | 17.96 <sup>±0.7</sup>       | 0.38 <sup>±0.01</sup>       | 7.10 <sup>±0.0</sup>        | 16.15 <sup>±0.19</sup>       | 0.04 <sup>±0.02</sup>       |

be attributed to excluding a corn mask, leading to mixed landuse/cover in the patches and noisy signals. Applying a corn mask directly on the input creates spatial gaps due to non-contiguous corn farms.

#### 4.3. Benefiting from improved spatial resolution

Using Landsat presents technical overhead, mainly in storage and computational requirements. For instance, a single-date Landsat observation for a county is 47 megabytes, approximately 280 times larger than MODIS. By comparing Landsat and MODIS (Scenario-2) under this setup, the highest performance is often seen with MODIS inputs except for the RF baseline. Similar to [9] there was no consistent value in using an improved resolution data for county-level corn yield prediction in large farming systems. Although Landsat offers a better level of spatial detail, it comes at the expense of reduced temporal resolution. Here, critical phenological states which may hold information necessary to estimate yield amount may be missed. Conversely, at farm-level, [13] deduced that using Landsat improved wheat prediction suggesting the appropriateness of higher spatial resolution for small spatial units.

#### 4.4. In-season prediction

Tables 3 and 4 present the RMSE and R<sup>2</sup> metrics at the end of selected months using LSTM for Scenario-1. The year 2012 can be predicted well in advance (end of July), but four units can improve the performance for 2020 at the end of August.

This improvement is achieved approximately 2-4 weeks before harvest, depending on the state.

**Table 3.** Comparing in-season RMSE using LSTM

|      | June  | July  | August | September    |
|------|-------|-------|--------|--------------|
| 2012 | 62.61 | 27.80 | 28.79  | <b>27.75</b> |
| 2020 | 18.18 | 17.58 | 13.34  | <b>12.78</b> |

**Table 4.** Comparing in-season R<sup>2</sup> using LSTM

|      | June  | July        | August | September   |
|------|-------|-------------|--------|-------------|
| 2012 | -1.40 | <b>0.53</b> | 0.49   | <b>0.53</b> |
| 2020 | 0.31  | 0.36        | 0.63   | <b>0.66</b> |

## 5. CONCLUSION

This study examined the benefits of improved spatial resolution for corn yield prediction by comparing MODIS and Landsat surface reflectance, complemented with weather variables. Our results suggest that spatially high, temporally low-resolution data offers no advantages for county-level yield assessments, while low-spatial, high temporal resolution data yields beneficial outcomes.

Regarding spatial-temporal feature extraction, the poor performance of the CNN-LSTM can be attributed to mixed or noisy land use in sampled patches, requiring an explanation layer to gain insights into the model's attention. In the case of Landsat, persistent cloud conditions can render an entire sequence of patches unusable. Such influence is minimal with

