

A deep learning approach for deriving winter wheat phenology from optical and SAR time series at field level

Felix Lobert^{a,b,*}, Johannes Löw^c, Marcel Schwieder^{a,b}, Alexander Gocht^a, Michael Schlund^d, Patrick Hostert^{b,e}, Stefan Erasmí^a

^a Thünen Earth Observation (ThEO), Thünen Institute of Farm Economics, Bundesallee 63, 38116 Braunschweig, Germany

^b Earth Observation Lab, Geography Department, Humboldt-Universität zu Berlin, Unter den Linden 6, 10099 Berlin, Germany

^c Department of Geoecology, Institute of Geosciences and Geography, University of Halle-Wittenberg, Von-Seckendorff-Platz 4, 06120 Halle (Saale), Germany

^d Faculty of Geo-Information Science and Earth Observation (ITC), University of Twente, Hengelosestraat 99, 7514 AE Enschede, the Netherlands

^e Integrative Research Institute of Transformations of Human-Environment Systems (IRI THESys), Humboldt-Universität zu Berlin, Unter den Linden 6, 10099 Berlin, Germany

ARTICLE INFO

Edited by Dr. Marie Weiss

Keywords:

Agriculture
Crop monitoring
Convolutional neural networks
U-net
Multisensor
Data fusion

ABSTRACT

Information on crop phenology is essential when aiming to better understand the impacts of climate and climate change, management practices, and environmental conditions on agricultural production. Today's novel optical and radar satellite data with increasing spatial and temporal resolution provide great opportunities to derive such information. However, so far, we largely lack methods that leverage this data to provide detailed information on crop phenology at the field level. We here propose a method based on dense time series from Sentinel-1, Sentinel-2, and Landsat 8 to detect the start of seven phenological stages of winter wheat from seeding to harvest. We built different feature sets from these input data and compared their performance for training a one-dimensional temporal U-Net. The model was evaluated using a comprehensive reference data set from a national phenology network covering 16,000 field observations from 2017 to 2020 for winter wheat in Germany and compared against a baseline set by a Random Forest model.

Our results show that optical and radar data are differently well suited for the detection of the different stages due to their unique characteristics in signal processing. The combination of both data types showed the best results with 50.1% to 65.5% of phenological stages being predicted with an absolute error of less than six days. Especially late stages can be predicted well with, e.g., a coefficient of determination (R^2) between 0.51 and 0.62 for harvest, while earlier stages like stem elongation remain a challenge (R^2 between 0.06 and 0.28). Moreover, our results indicate that meteorological data have comparatively low explanatory potential for fine-scale phenological developments of winter wheat.

Overall, our results demonstrate the potential of dense satellite image time series from Sentinel and Landsat sensor constellations in combination with the versatility of deep learning models for determining phenological timing.

1. Introduction

Phenology refers to the study of periodic events in the life cycle of organisms, which are mainly triggered and controlled by environmental factors (Lieth, 1974; Morissette et al., 2009). When monitoring plants and in particular crops, information on seasonal phenology allows understanding a crop's metabolic cycle, its response to meteorological drivers such as temperature and humidity, and its buildup of biomass, among others (Richardson et al., 2013). Crop phenology is hence a valuable

input for numerous agricultural monitoring tasks, including the assessment of management practices and yield estimation. Furthermore, phenological information is a reliable indicator for climate change impact analysis and of high interest in fields like ecology and global change biology (Ma et al., 2022; Menzel, 2002; Menzel et al., 2006).

Meaningful and large-scale analyses of phenological patterns require a large amount of in-situ data, whose field-based collection is hardly feasible. Gerstmann et al. (2016) demonstrated the potential of meteorological data to map general phenological patterns for several crops

* Corresponding author at: Thünen Earth Observation (ThEO), Thünen Institute of Farm Economics, Bundesallee 63, 38116 Braunschweig, Germany.
E-mail address: felix.lobert@thuenen.de (F. Lobert).

based on the well-known relations between temperature, precipitation, and plant development. In comparison to meteorological data, Earth Observation (EO) satellites directly capture the condition of vegetation at the field level and thus provide proximate information on plant development. These temporal signals of vegetation development revealed by satellite sensors were defined as land surface phenology (LSP; De Beurs and Henebry, 2004). Such satellite-based observations of LSP allow to infer phenological changes of crops on the ground and derive phenological information. The field of satellite-based phenology research has been around for a long time, yet new methods like Deep Learning (DL) and possibilities of sensor fusion represent potentials that have not yet been fully exploited (Katal et al., 2022; Pipia et al., 2022). Therefore, it is of great interest for the EO research community to further investigate these potentials and contribute to spatially and temporally improved proxies from satellite data analyses for identifying crop phenological stages.

Studies focusing on phenology analysis based on satellite data usually aim at identifying specific points in remote sensing time series that represent key events of the crop's life cycle, such as the Start Of Season (SOS) or End Of Season (EOS; Zeng et al., 2020). This is often achieved by defining thresholds for Vegetation Indices (VI) that can either be static or dynamic (e.g., Bolton et al., 2020; Meroni et al., 2021). Another way is to calculate derivatives from satellite data time series that can be used to identify breakpoints, turning points, or other significant changes in the trend of the time series, which are mainly inspired by mathematical curve descriptors (Harfenmeister et al., 2021; Kowalski et al., 2020; Schlund and Erasmi, 2020). The products resulting from these methods can be understood as phenological metrics that provide estimates for the general progression during plants' life cycles and are, therefore, of great interest for various applications. However, these phenological metrics are rather mathematical descriptors of VI curves and are not necessarily linked to the sharply defined phenological events we can measure in the field, like stem elongation or heading of wheat (Zhang et al., 2017). Applications, such as biophysical plant growth or yield models, that need detailed information about individual phenological stages, therefore, require new methods to provide such input.

Traditionally, optical imagery has been used as predominant data source for deriving phenology information from remote sensing. Time series of VIs and raw band measurements show characteristic patterns that can be attributed to changes in the plant as it progresses through the various phenological stages, such as the fraction of ground cover, chlorophyll content, and color. However, optical imagery usually comes with the issue of data gaps in time series introduced by clouds and cloud shadows. Data gaps hamper the detection of changes in a crop's temporal signature. Thus, there has been a trend in phenological analyses towards using synthetic aperture radar (SAR) data, specifically since the advent of operational Sentinel-1 data (e.g., Löw et al., 2021; McNairn et al., 2018; Nasrallah et al., 2019; Schlund and Erasmi, 2020). Derived features, like the backscatter coefficient, are sensitive to surface roughness and the dielectric constant. These properties depend on vegetation structure, leaf angles, vegetation cover, and water content, which change during the phenological development of crops. Mainly during the first months after seeding, soil roughness and moisture potentially influence the SAR signal.

One of the advantages of SAR data against optical data time series, is that they are usually not impeded by data gaps due to cloud cover. However, speckle noise, precipitation-induced soil moisture changes, and different acquisition geometries are common challenges when working with SAR data and limit time series quality. Consequently, combining spectral and textural/structural information derived from both optical and SAR systems can help mitigate the weaknesses of each data type and create synergies instead (Meroni et al., 2021; Pipia et al., 2022).

The suitability of optical and SAR data for phenological analyses was already investigated and compared in several studies (d'Andrimont et al., 2020; Fieuzal et al., 2013; Meroni et al., 2021; Nasrallah et al.,

2019; Veloso et al., 2017). Most of them agree on the potential arising from the joint use of data from both sensor types. However, studies presenting methods that make use of this combination were only introduced recently by Mercier et al. (2020) and Yeasin et al. (2022). Both reported improvements over single-sensor models, supporting the assumption of data complementarity for phenological analyses. However, Mercier et al. (2020) and Yeasin et al. (2022) were based on a limited number of observations, which hampers accurate inferences about the timing of actual stage transitions.

Nowadays, we are faced with a wealth of data from various Earth observation missions. However, we still need advanced methods that can appropriately exploit their potential to estimate phenological information on arable crops. DL has been shown to be a suitable tool for the combined exploitation of multivariate time series from heterogeneous data sources (Holtgrave et al., 2023; Lobert et al., 2021), while the potential of DL for multi-sensor phenological analyses is generally under-studied (Katal et al., 2022). We here consequently address this research gap by utilizing a supervised one-dimensional DL model that is inspired by phenology-like problems in medical time series applications (Jimenez-Perez et al., 2019; Perslev et al., 2019). We exploited data from Sentinel-1 (S1), Sentinel-2 (S2), and Landsat 8 (L8) together with meteorological data and a comprehensive data set on phenological field observations provided by the German Weather Service (DWD). Being the most widely grown crop in Germany, we focused on winter wheat (Federal Statistical Office, 2022). We compared around 16,000 phenology observations to nearby field-level remote sensing time series for winter wheat between 2017 and 2020. The model was then trained to predict the start of seven different phenological stages at field level and compared against a baseline provided by a Random Forest (RF) classifier (Breiman, 2001).

The presented approach contributes to innovation in the field of crop phenology estimation in two respects: first, the chosen architecture represents an "all-in-all-out" approach, i.e., time series of different features are simultaneously fed into the model that predicts the entry data of multiple phenological stages at once. This extends the current state-of-the-art that mostly builds on separate rule sets or features for different stages (Zeng et al., 2020). Second, the model training enables us to directly search for relevant patterns in the time series instead of defining the key points ex-ante and matching them to field observations afterward.

We hence aimed to answer three research questions:

1. What is the performance of the proposed one-dimensional DL model to predict the start of phenological stages for winter wheat at field level based on different sets of input features and against the baseline model?
2. How does the performance differ for the individual stages?
3. How do our estimates of the start of phenological stages compare to spatiotemporal patterns of the ground observations across Germany?

2. Study area and data

2.1. Study area and reference data

For our study, we used reference data provided by DWD. Around 1200 trained volunteers located across Germany observe the phenology of nearby plants (Kaspar et al., 2015). The volunteers choose one field for each crop within a distance of 2 km (up to 5 km in exception) from their reported base location (Fig. 1; Deutscher Wetterdienst, 2015). An assignment to a specific field, however, is not provided. We selected the observations from around 700 volunteers who surveyed the start (reached on 50% of the field) of seven different phenological stages for winter wheat, always on the same field (Fig. 1; DWD, 2022a). The observations begin with the seeding of the winter wheat, followed by the start of leaf development, stem elongation, heading, milk ripeness, yellow ripeness, and lastly harvest. Our study covered four vegetation

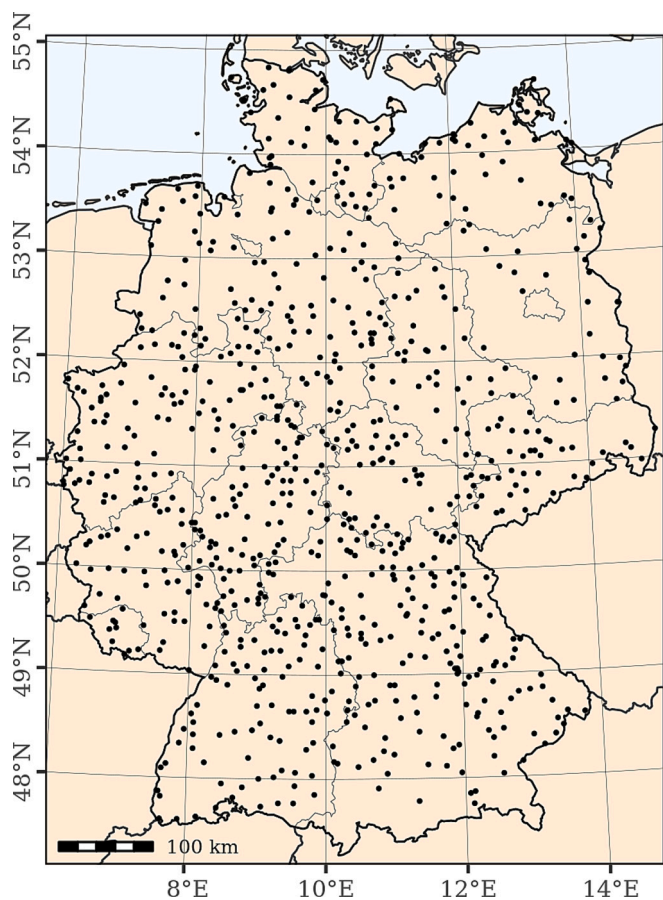


Fig. 1. Locations of the phenological observations in Germany (DWD, 2015, 2022a).

periods from the seeding of winter wheat in autumn 2016 to the harvest in late summer 2020. The combination of 700 observation stations, seven stages, and four observed vegetation cycles results in over 16,000 observations.

The locations of the observations cover the full gradient of climate and topographic characteristics across Germany, from the Alpine foreland in the South, over regions with a continental climate in the East to a maritime climate in the West and the Northern German lowlands. The observation period (2016 to 2020) covers heterogeneous meteorological conditions. While the year 2017 experienced average amounts of precipitation and temperatures in Germany, 2018 was exceptionally dry and hot (Fig. 2). Subsequent years 2019 and 2020 were also characterized by low to average moisture conditions, which prevented recharge of groundwater storage.

The start of the phenological stages during the four studied years reflects the described climate conditions. (Fig. 3). The timing of seeding and leaf development throughout the study period does not show large deviations or outliers. The start of stem elongation occurred later in 2018, which could be related to cold conditions at the beginning of 2018. As of April 2018, very hot and dry conditions can be observed as well as a much earlier start of the heading to harvest stages compared to the other years.

2.2. Field boundaries

We used a German-wide crop type map (CTM), which was produced by Blickensdörfer et al. (2022) based on S1, S2, and L8 data for identifying the main crop types in Germany at 10 m spatial resolution. For each observer location, we extracted all winter wheat pixels from the CTM of the respective year within a surrounding of 5 km. Adjacent pixels were then clustered and combined into individual fields. To enhance the quality of the field boundaries, we utilized a two-step buffering approach. First, we applied an inward buffer of 70 m to each boundary. This was then followed by an outward buffer of 40 m. This procedure imitates a morphological opening operation and removes erroneous connections between multiple fields. Using a higher value for inward

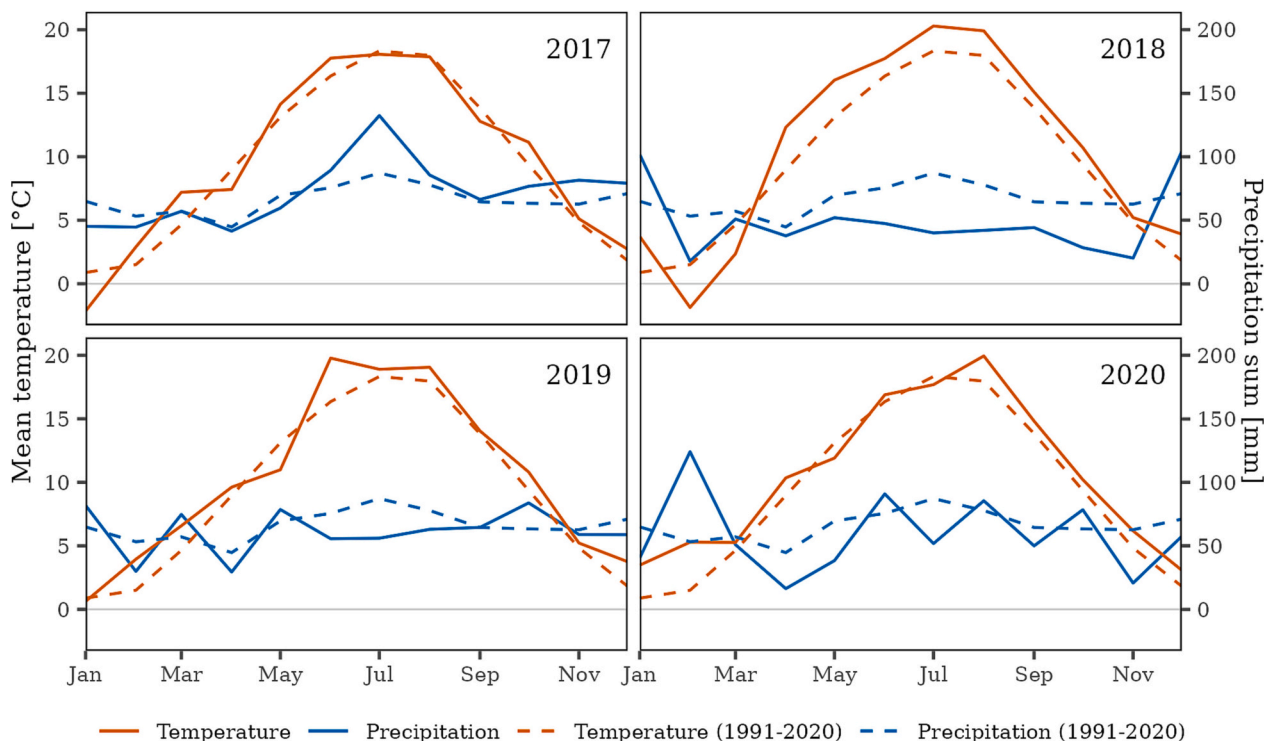


Fig. 2. Monthly mean records of temperature and precipitation sums in Germany during the studied years and long-term average (1991–2020) (Source: DWD).

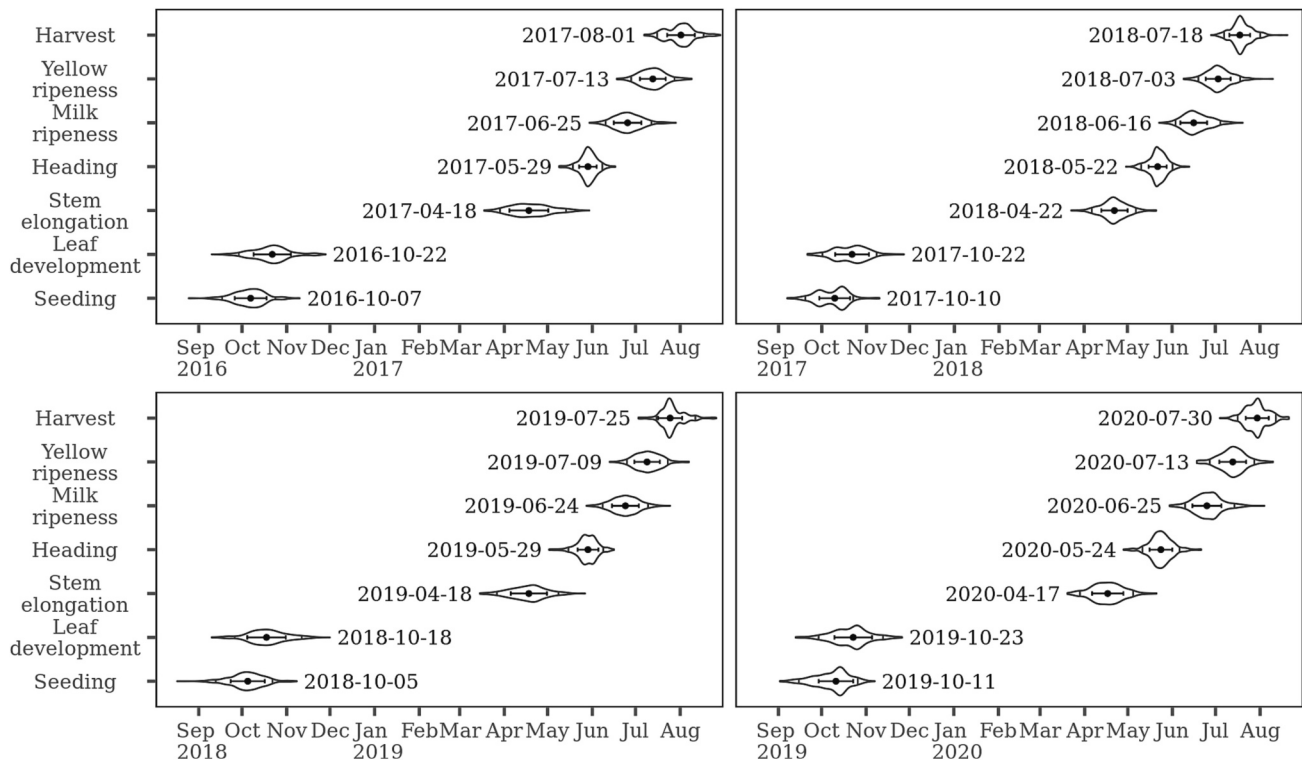


Fig. 3. Temporal distribution of the phenological observations for winter wheat in the studied growing seasons. Points represent the median (annotated date), error bars show ± one standard deviation. Vertical bars in the violin plots show the 5th and 95th percentiles.

buffering mitigates edge effects along the field boundaries. Finally, fields smaller than 2 ha were discarded, to exclude excessively small fields from the training process. In addition, we decided to limit our analysis to the 10 closest fields to the observer's position to reduce the bias by high variations in the number of winter wheat fields and not weigh distant fields too much, if there are already enough fields close to the observer location (Fig. 4). This procedure resulted in about 22,000 field boundaries for winter wheat that were linked to the phenological observations during one growing season.

2.3. Remote sensing imagery

2.3.1. Sentinel-1

We used the gamma naught (γ^0) backscatter coefficient from the S1A and S1B constellation as SAR-based input. S1 acquires data in the C-band (5.4 GHz, 5.5 cm), with dual polarization mainly in VV (vertical transmit and vertical receive) and VH (vertical transmit and horizontal receive). Standard acquisitions are in interferometric wide swath (IW) mode, which covers a swath of about 250 km (Torres et al., 2012). We used the Ground Range Detected (GRD) IW product.

Since the launch of S1B in 2016 and until its unexpected failure at the

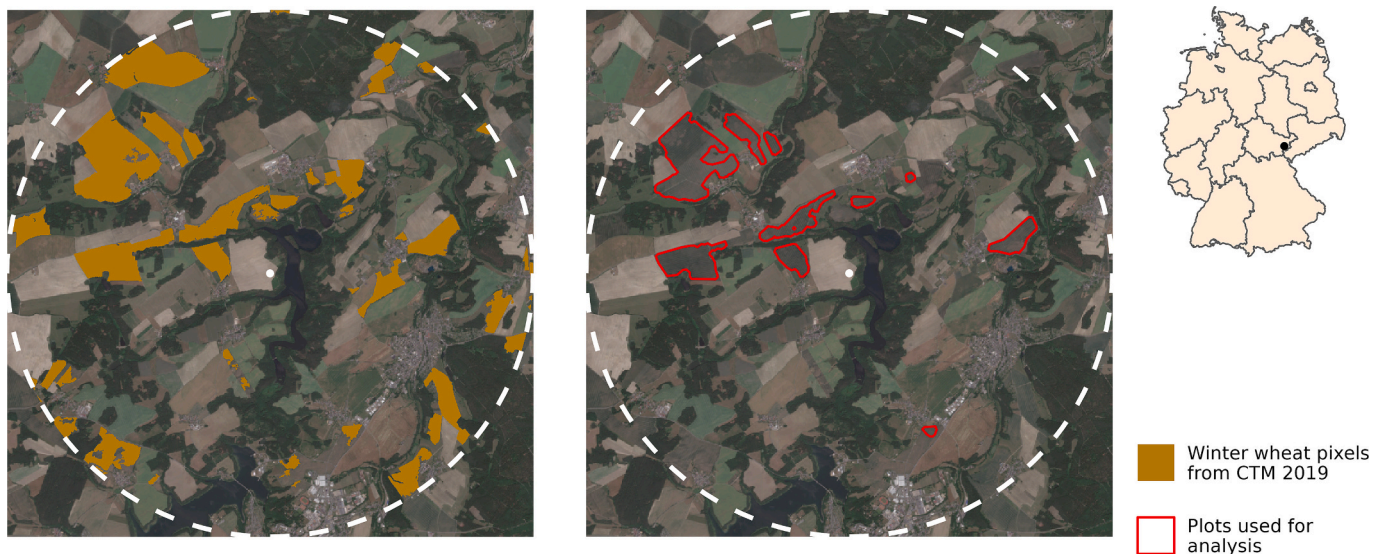


Fig. 4. Example of a reported observer location from the reference data (white point), extracted winter wheat pixels within 5 km distance (white dashed circle), and resulting field boundaries in 2019 that were used for further analysis (right). Background image: monthly RGB-composite from Sentinel-2 for June 2019.

end of 2021, the S1 constellation acquires data at a 6-day interval. We used all available data from both sensors and across all orbits for Germany during our study period. This resulted in 18,203 S1 scenes from August 2016 to October 2020. S1 accordingly delivered an observation every 1.8 days on average, depending on orbit overlap across Germany. We accessed the S1 data through the Copernicus Data and Exploitation Platform - Germany (CODE-DE; Benz et al., 2020). The pre-processing was carried out using the Sentinel Application Platform (SNAP) and the R package *rcodede* (Lobert, 2022).

The γ^0 backscatter coefficient was processed by first applying border and thermal noise removal to the S1 GRD scenes. This was followed by calibration and radiometric flattening of the data to obtain the γ^0 backscatter coefficient in VV and VH polarization in dB. Gamma naught represents the ratio between the incident power and the scattered power for a reference area that is perpendicular to the line of sight from the sensor to an ellipsoidal model of the ground surface (Small, 2011). The imagery was terrain corrected using the Shuttle Radar Topographic Mission (SRTM) 1 arc-second global digital elevation model (DEM; Farr et al., 2007), and resampled to 10 m spatial resolution.

We then calculated the backscatter cross-ratio (CR) to exploit the information content of the backscattered signal in both polarizations

$$CR = \gamma_{VH}^0 [dB] - \gamma_{VV}^0 [dB] \quad (1)$$

which is strongly affected by structural changes in crops like winter cereals (Holtgrave et al., 2020; Nasrallah et al., 2019; Vreugdenhil et al., 2018). Moreover, Schlund and Erasmi (2020) reported that the CR produces a relatively stable signal in dense time series over longer periods over agricultural areas since both polarizations react similarly to terrain and soil properties which reduces the impact of these factors on the CR signal. Meroni et al. (2021) have shown that this also allows for the combined use of multiple orbits and acquisition directions, enabling the analysis of time series consisting of up to daily observations in areas of orbit overlaps.

2.3.2. Sentinel-2 & Landsat 8

We obtained L8 as Level-L1TP and S2 as Level-1C data. We used all available scenes that cover Germany during our study period with a cloud coverage of less than 75% and corrected all data for radiometric and geometric effects using the Level 2 processing system in FORCE

(Frantz, 2019). Clouds and cloud shadows were masked out using the improved Fmask algorithm (Frantz et al., 2018; Zhu et al., 2015; Zhu and Woodcock, 2012).

We applied a spectral adjustment between S2 and L8 according to Scheffler et al. (2020). Spectral harmonization uses S2A as reference and adjusts the spectral response of S2B and L8 to S2A, including a prediction of missing Sentinel bands for L8. Bands for atmospheric correction, as well as panchromatic and thermal bands of the optical sensors were not further considered. The Enhanced Vegetation Index (EVI) was calculated to complement the original spectral bands (Huete et al., 2002).

We organized the data in a tiled and reprojected data cube. We resampled all imagery to 20 m spatial resolution using nearest neighbor resampling. We ended up with an average clear sky observation (CSO) for our fields approximately every 7 days. Spatial and temporal patterns emerging from orbit overlaps or sensor availability are visualized in Fig. 5 and Fig. 6.

2.4. Meteorological data

We used daily mean temperature measurements (°C) from 625 weather stations across Germany (DWD, 2022b). Precipitation data was acquired from the German weather radar network RADOLAN, which provides area-wide rainfall estimates with a temporal resolution of 5 min and a spatial resolution of about 1 km (DWD, 2022c). Data from DWD were accessed and preprocessed using the *rdwd* R package (Boessenkool, 2021).

3. Methods

The analysis concept of our study relies on the association of phenological observations, crop type information, and various remote sensing and environmental time series to train a supervised classification model and conduct an analysis of feature importance. The detailed procedure is depicted in Fig. 7 and the following sections. We carried out all processing steps using R (R Core Team, 2022).

3.1. Time series preprocessing & labeling

Since the analysis centers on the field level, we transformed the areal

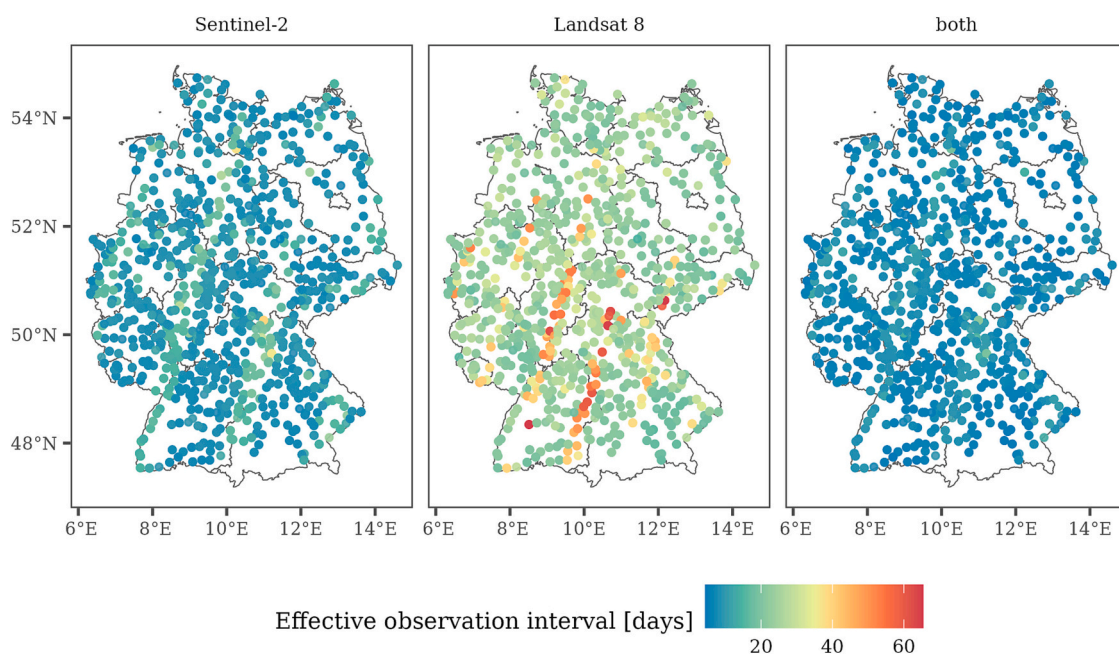


Fig. 5. Average interval between two clear sky observations (CSO) for the observer locations during the years 2017–2020.

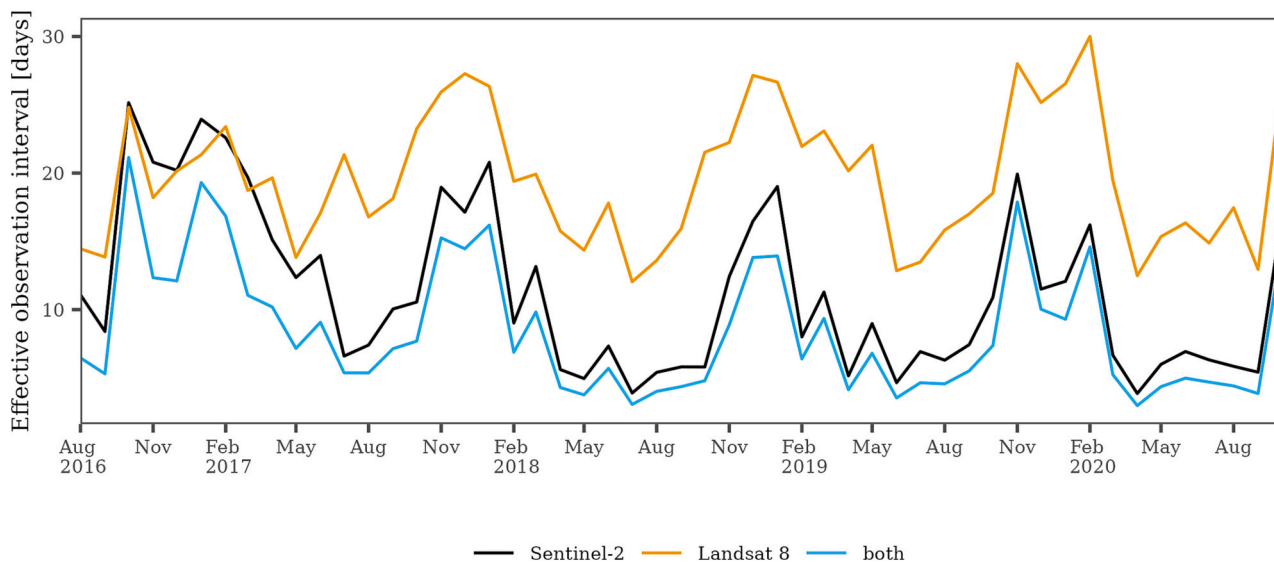


Fig. 6. Temporal distribution of the interval between two clear sky observations (CSO) averaged per month.

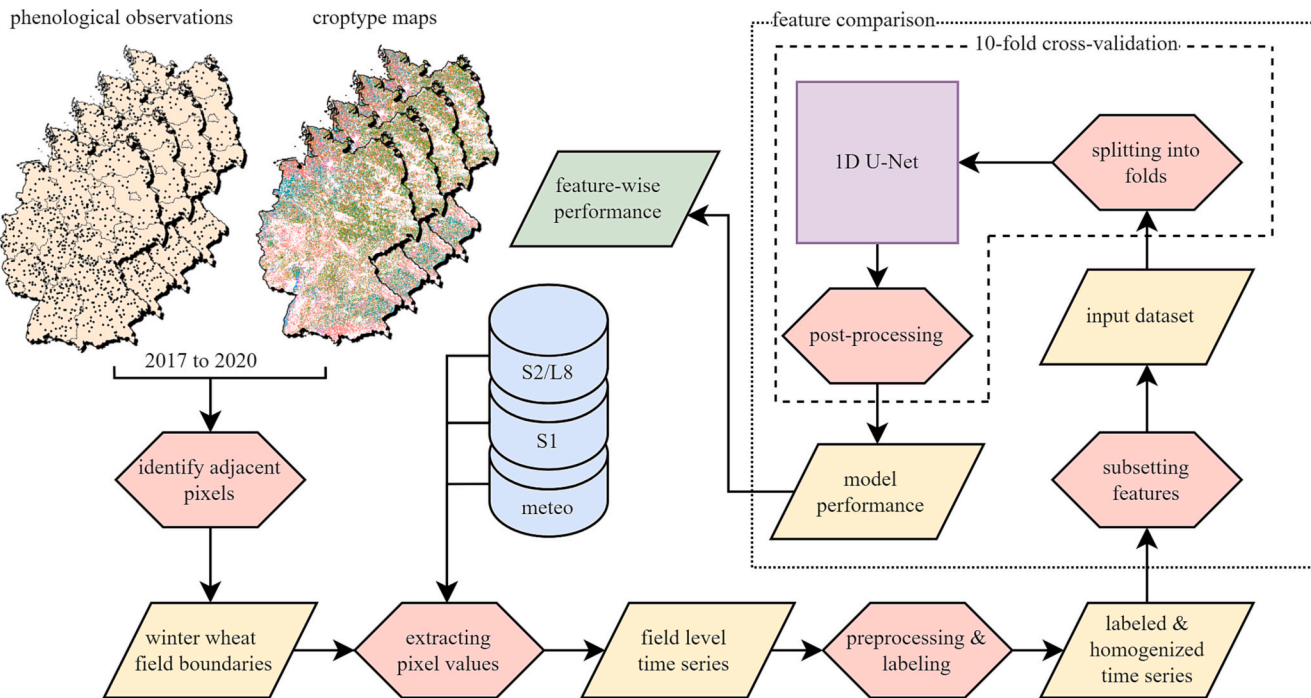


Fig. 7. Workflow of the method proposed in this study.

information from the remote sensing imagery to one-dimensional time series per field. This was realized by summarizing the pixel values for each field, date, and input feature using the field boundaries. We chose the median to account for outliers. We performed this for the gridded input data, including the optical and SAR imagery as well as the RADOLAN precipitation data. We interpolated the temperature measurements from the ten nearest DWD weather stations for each field using the inverted distance weighting method (IDW) and inverse distance power set to 0.5. We acquired time series for each field starting in August before sowing and ending at the end of November of the following year.

We then applied locally estimated scatterplot smoothing (loess) to account for undesirable noise and artifacts in each time series (Cleveland et al., 1992). A *span* parameter of 0.3 was visually assessed to yield the

best trade-off between preserving enough information and suppressing noise. As our chosen model architecture required equidistant time steps, we considered a three-day interpolation interval to be apt for phenology monitoring. We realized this through linear interpolation of the optical, SAR, and temperature data, while the precipitation data were summed up for the last three days preceding every time step. We finally normalized values per feature, field, and growing season by subtracting the mean and dividing by the standard deviation. This improved the comparability across different fields and years and ensured that all of the features were in the same value range to ease the learning process during model training (Bishop, 1995). A composition of exemplarily pre-processed time series for a winter wheat field from the seeding in 2017 to the harvest in 2018 is shown in Fig. 8.

We finally added labels to each 3-day time step of our time series. For

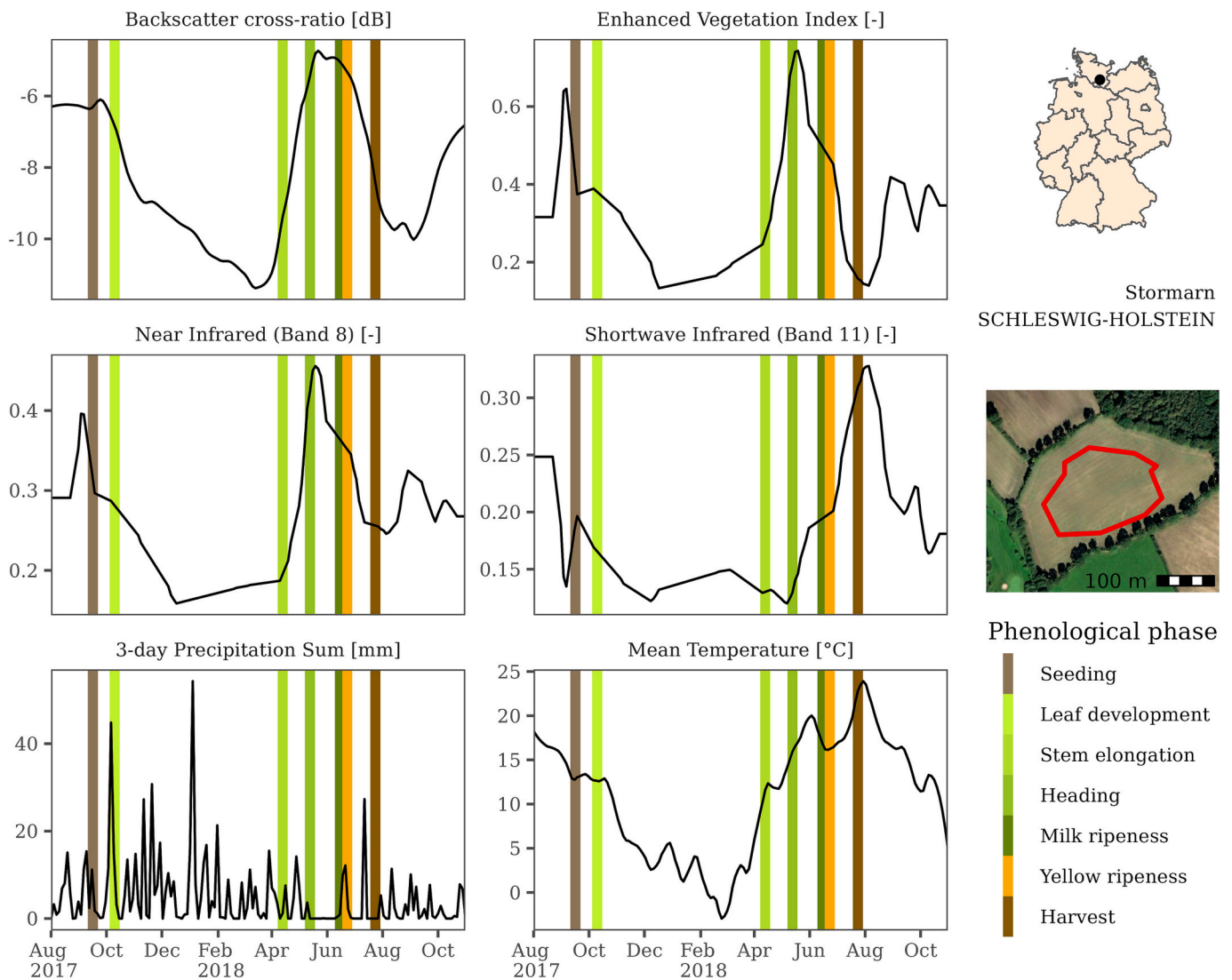


Fig. 8. Time series of different exemplary features for a winter wheat field from seeding in 2017 to harvest in 2018. Vertical lines show the observed start of the phenological stages.

each time series, we identified the closest time steps to each reference stage from the corresponding DWD observation and labeled the time steps accordingly. We extended these labels by ± 3 days, resulting in three labeled time steps for each stage. All other time steps were labeled as background class, where none of the recorded stage changes took place.

3.2. Deep learning model

Convolutional Neural Networks (CNNs) are a commonly used model architecture in remote sensing of vegetation (Kattenborn et al., 2021). CNNs usually consist of multiple convolutional layers that can be connected in different ways. Due to the nature of the convolution process, these layers are ideal for detecting changes in sequential data. For this reason, CNNs are prominent, e.g., for the detection of boundaries in two-dimensional data structures like images. Although CNNs are mainly used in a two-dimensional design, convolutions can also be used to analyze one-dimensional data, such as time series. This type of use was already demonstrated to be powerful for classification and event detection tasks when dealing with pixel- or field-based time series of satellite data (Lobert et al., 2021; Pelletier et al., 2019).

Ronneberger et al. (2015) proposed the U-Net architecture, which is based on multiple interconnected convolutional layers that analyze data

in different aggregation levels. Jimenez-Perez et al. (2019) and Perslev et al. (2019) adapted the U-Net architecture to one-dimensional data, transferring the U-Net's ability to delineate object borders in images to delineate processes and events in time series. They used their adapted architecture to detect and delineate cardiac illnesses from electrocardiograms (ECG) and sleep stages from electroencephalograms (EEG).

3.2.1. Implementation

Inspired by these developments, we implemented our own one-dimensional U-Net architecture to predict the start of different phenological stages of winter wheat at the field level. Starting from the classic architecture of the U-Net by Ronneberger et al. (2015), our first major change was to adapt the input layer to read our time series of 152 3-day time steps and multiple features. This corresponds to the time series length over the extended winter wheat growing season and the different input features derived from remote sensing and meteorological data. As a second major modification, we replaced every second convolutional layer in the down- and upsampling path of the U-Net with a Long Short-Term Memory layer (LSTM; Hochreiter and Schmidhuber, 1997). These recurrent layers allow for more powerful exploitation of the temporal domain of data that is beyond the length of the convolutional filter kernels. The final architecture of our model including the filter numbers of the convolutional layers and the amount of LSTM cells is depicted in

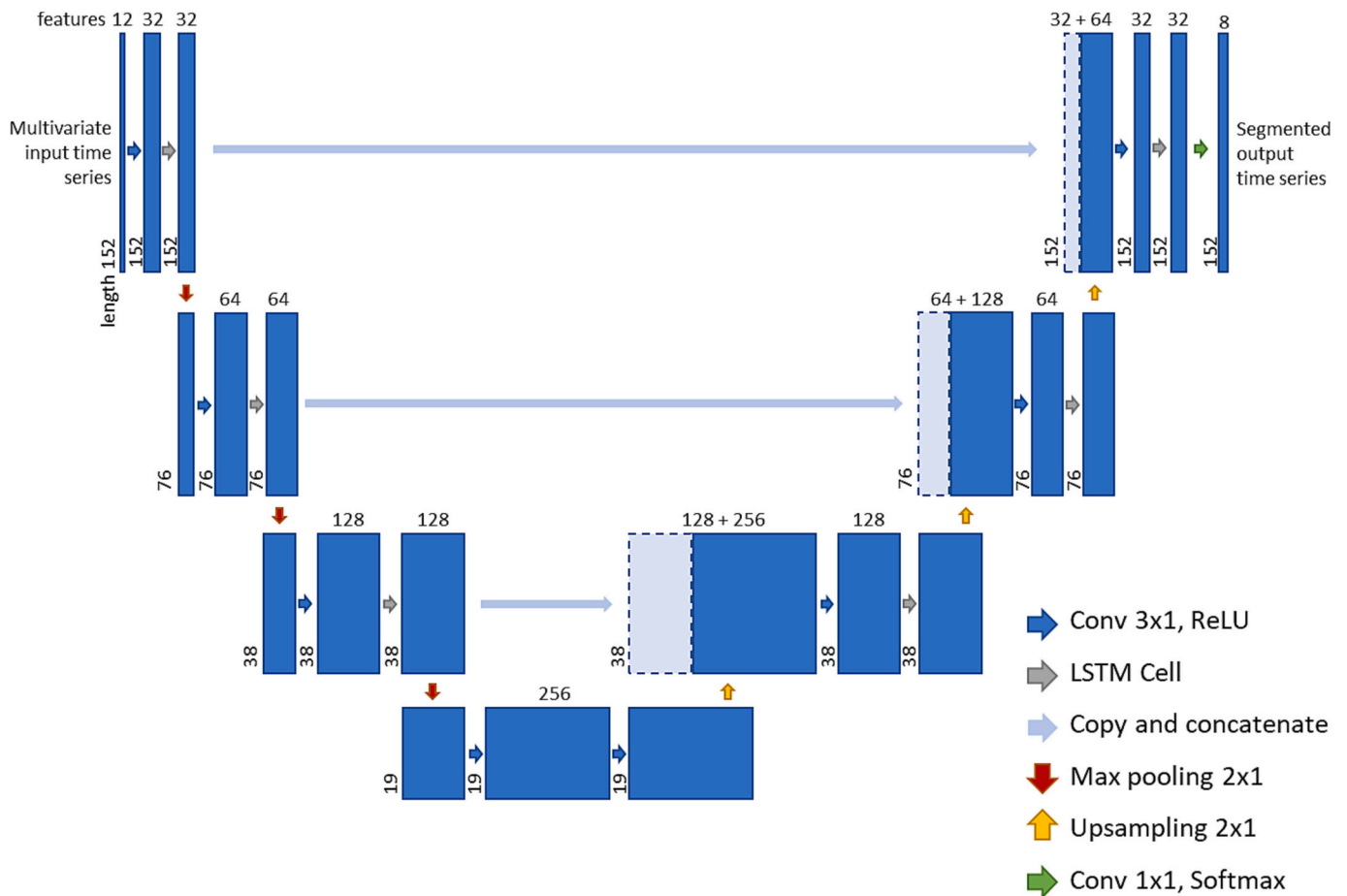


Fig. 9. Schematic architecture of the model used in our study adapted from the initial U-Net design by Ronneberger et al. (2015). Here, an example with 12 input features and eight output classes is shown (7 phenological stages plus background class). In contrast to the U-Net model, every second convolutional layer in the up- and downscaling layers is replaced by an LSTM layer.

Fig. 9. The final output of our model is a time series of the same length as the input for the respective phenological stage. The output values provide the probability of each time step to be the start of the respective stage. We used the rectified linear unit (ReLU) activation function for convolutional layers and hyperbolic tangent (tanh) activation for LSTM layers to speed up the training process on a graphical processing unit (GPU) and activated the model output using softmax. We implemented our model using Keras (Chollet, 2015) with TensorFlow (Abadi et al., 2016) as backend on the R interface to Keras (Allaire and Chollet, 2021).

3.2.2. Training

We used three independent data sets, i.e., training, validation, and test data for building our model. We used the training data to train the model and perform the backpropagation. After each training epoch, the model was applied to the validation data to provide insights into the generalization ability and to allow for the adaptation of optimization parameters during training. In a subsequent step, the model was evaluated using the test data that were previously unseen by the model.

We used categorical cross-entropy as loss function, to respect not only the correctness of the predicted classes but also the certainty of the predictions. For the calculation of the loss function, we used temporal sampling weights to weigh those errors higher that were closer to phenological stage transitions. We employed Adam (Kingma and Ba, 2015) as optimization algorithm with an initial learning rate of $1e^{-3}$ and a batch size of 2^8 . We performed 200 training epochs but applied an early stopping mechanism to end the training when the loss function did not decrease for 50 epochs. All other parameters remained as Keras defaults.

3.3. Evaluation

Our study aimed to evaluate the performance of the proposed model for detecting phenological stages given different sets of input features. We, therefore, defined five different feature sets that were tested during our validation (Table 1).

Table 1
Feature sets that were tested in this study.

feature set	input features	number of features
SAR	γ_0 backscatter coefficient VV [dB] γ_0 backscatter coefficient VH [dB] backscatter cross-ratio [dB]	3
optical	blue (496.6 nm) [-] green (560.0 nm) [-] red (664.5 nm) [-] red edge 1 (703.9 nm) [-] red edge 2 (740.2 nm) [-] red edge 3 (782.5 nm) [-] near-infrared (835.1 nm) [-] shortwave infrared 1 (1613.7 nm) [-] shortwave infrared 2 (2202.4 nm) [-] enhanced vegetation index [-]	10
meteorological	precipitation sum [mm] mean temperature [°C]	2
SAR & optical	SAR features optical features	13
all	SAR features optical features meteorological features	15

To get an estimate of our overall model performance, we decided to conduct our evaluation based on 10-fold cross-validation (CV). We randomly sampled our input data into 10 equally sized folds, thereby ensuring that all time series belonging to the same phenological observation ended up in the same fold. We went for random CV since spatial CV approaches can lead to overly pessimistic accuracy estimates. This is because whole geographic regions and with this, environmental conditions and also regionalized agricultural management practices are left out during the training process in each cycle of the CV. This was shown by Wadoux et al. (2021), who observed no improvement in spatial CV over random strategies in their comparative study. Furthermore, random CV is less of an issue if the model is not intended to extrapolate but to be applied within the environmental range of the training data (Kattenborn et al., 2022). Here, we used each of the folds as test data for one training cycle, while the remaining folds became the training data (80%) and validation data (20%).

We transformed the predicted probabilities for each field, time step, and phenological stage into discrete predicted dates for their start before finally evaluating the model results. This was realized by first searching for the time step with maximum probability for each phenological stage. We then selected the five preceding and following time steps and calculated the mean of the dates, weighted by their probabilities. This procedure allowed making predictions with a finer temporal resolution than the temporal interval of our time series. The output was rounded to a (full) day of year (DOY) and finally formed our discrete predictions. An example of the transformation from probabilities to discrete predictions is shown in Fig. 10.

DWD field measurements are conducted on the same winter wheat field throughout the growing season but the provided dataset lacks assignment to a specific field (section 2.1). We identified up to ten candidate fields for each measurement during preprocessing (section 2.2). Obtaining individual predictions for each candidate field, leads to the need for a strategy to evaluate model performance. Averaging the predictions for all candidate fields eliminates the variance of the model predictions across different fields, potentially resulting in an overly

pessimistic performance estimate. To address this, we adopted the *minimum bias* approach proposed by Ye et al. (2022). We calculated the absolute error for each candidate field across all 7 phenological stages, identifying the field with the overall least bias. However, the *minimum bias* method may lead to overly optimistic results as the prediction selection is not completely independent of the reference data. Therefore, we evaluated our results using both approaches and discuss their differences. The first approach is referred to as *mean prediction*, while the second approach is referred to as *minimum bias prediction*.

For the validation, we first compared the performance of the different feature sets. We determined the accuracy of our predictions by considering them correct if they were made within a six-day window from the reference date. This measure, defined as prediction accuracy, represents the proportion of correctly predicted outcomes in relation to the total predictions made. We chose this time frame for technical reasons with respect to our time series interval and expected label noise. We compared the performance of the models trained with different features sets and performed McNemar's test to test for significance (McNemar, 1947). Based on the prediction accuracy, we identified the best-performing feature set, for which we then conducted a more in-depth analysis of the model performance. Calculating the mean absolute error (MAE) and the coefficient of determination (R^2) enabled us to compare the different phenological stages. We mapped spatial and temporal distributions of the predictions and analyzed emerging patterns. Furthermore, the temporal transferability of the model and differences between the years were assessed by performing an additional temporal cross-validation, where in each cycle one year was left out for training and instead used for testing.

To provide a baseline for comparison of the proposed DL model, we also tested a RF regression model for our task (Breiman, 2001). Since multidimensional input and output are not supported by RF, we flattened our input features and trained one model for each stage, using the DOY as the target variable. The R package caret (Kuhn, 2020) was used with the corresponding default parameters, and the same cross-validation scheme as for the U-Net. To ensure a more focused and

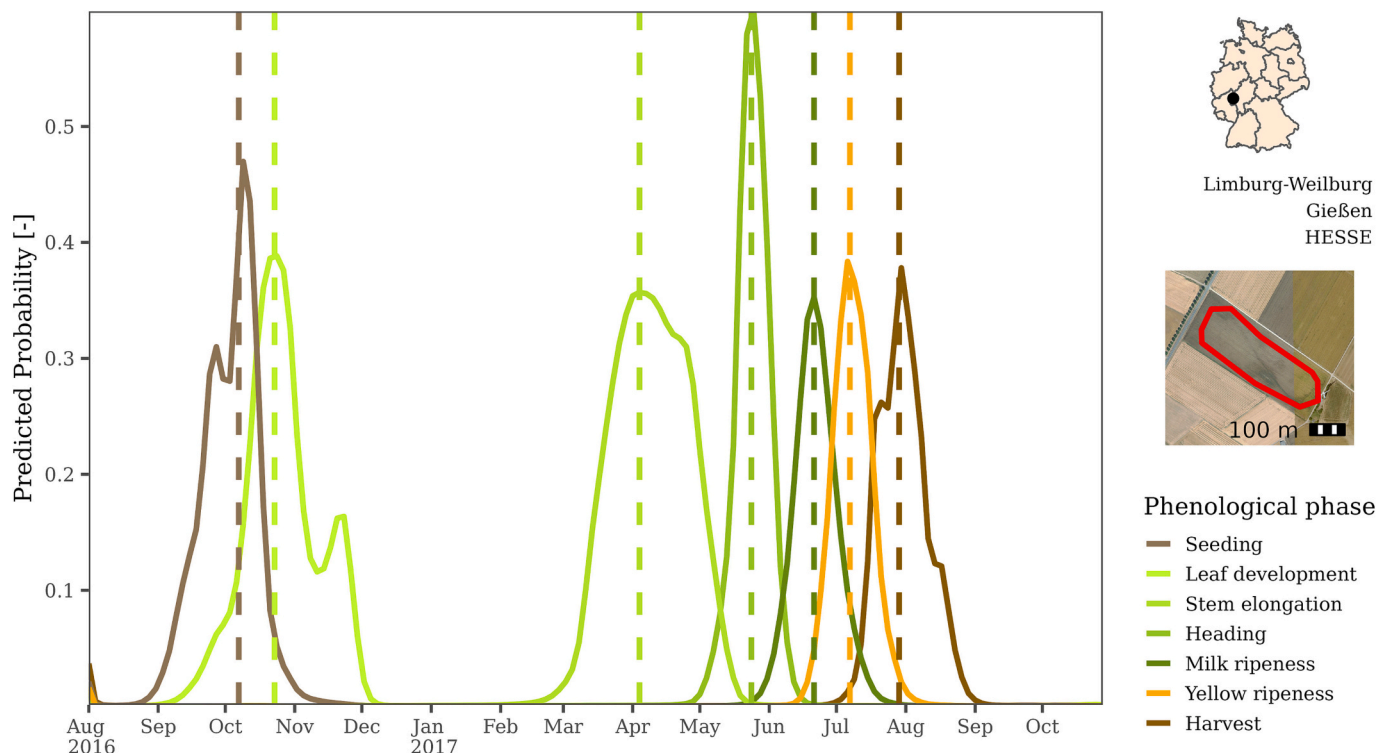


Fig. 10. Predicted probabilities for each time step to be the start of each phenological stage as predicted by the model for an exemplary field. The vertical dashed lines show the derived discrete predicted date.

efficient analysis, we limited our model comparisons to the best-performing feature set identified by the U-Net model, avoiding an excessive number of comparisons.

4. Results

4.1. Comparison of input features

The overall results of our feature set comparison are visualized in Fig. 11. On average, SAR and optical data performed similarly. Only a slightly higher prediction accuracy of 49.9% and 65.1% (*mean* and *minimum bias*) for SAR compared to 49.0% and 64.2% for the optical data was observed. However, we found differences between the individual stages. The highest differences occurred for the *minimum bias* predictions for seeding with 62.6% prediction accuracy for the SAR data set as compared to 56.4% for optical data, which is supported by a high level of significance according to McNemar's test. Heading also showed notably higher accuracies based on SAR data, especially for the *mean* predictions (SAR: 64.1%, optical: 59.6%) and significant differences. Yet, there were stages where optical data showed higher prediction accuracies, although not being significant. This was the case, especially for harvest with 58.2% and 73.7% (*mean* and *minimum bias*) for SAR compared to 60.7% and 75.6% for optical data. For the yellow ripeness stage, optical data only performed better when considering the *mean* prediction (SAR: 53.6%, optical: 55.7%) and for milk ripeness only when considering the *minimum bias* prediction (SAR: 62.6%, optical: 63.8%). Overall, radar data were performing better for the early phases, while optical data were ahead for the late phases.

Combining SAR and optical data did not show a clear improvement in the general model performance compared to solely using SAR data. On average, the prediction accuracy only increased from 49.9% (SAR) to 50.1% (SAR & optical) and 65.1% (SAR) to 65.5% (SAR & optical) for the *mean* and *minimum bias* predictions without significant differences. However, compared to optical data, the combination resulted in higher accuracies (*mean*: 49.0% to 50.1% and *minimum bias*: 64.2% to 65.5%) and significant differences in the predictions. Predictions for leaf development improved most, yet only for the *minimum bias* predictions with 55.8% for SAR and 58.2% for both features combined. The harvest stage also improved with 62.7% and 76.9% for the combination of SAR and optical data compared to optical data with only 60.7% and 75.6%, for *mean* and *minimum bias* predictions. Some stages, however, decreased in performance when both data sets were combined. This applies to seeding and heading, where SAR data alone performed better.

The meteorological feature set showed less explanatory power compared to the remote sensing-based data sources. This feature set yielded the lowest prediction accuracy both on average as well as for the individual stages and showed significant differences in all comparisons. This applies equally to the *mean* and *minimum bias* predictions. Adding meteorological data to the input features only improved the prediction accuracy for seeding (*mean*: 38.9% to 39.2% and *minimum bias*: 60.1% to 62.7%) and heading (*mean*: 59.7% to 63.8%). Yet, SAR data alone still performed better for seeding (*mean*: 41.4%) and heading (*mean*: 77.9%). Generally, the accuracies for the *minimum bias* predictions were higher than the *mean* predictions for all stages and feature sets. However, the difference was smaller for the meteorological feature set.

Based on this comparison, we identified the combination of SAR and

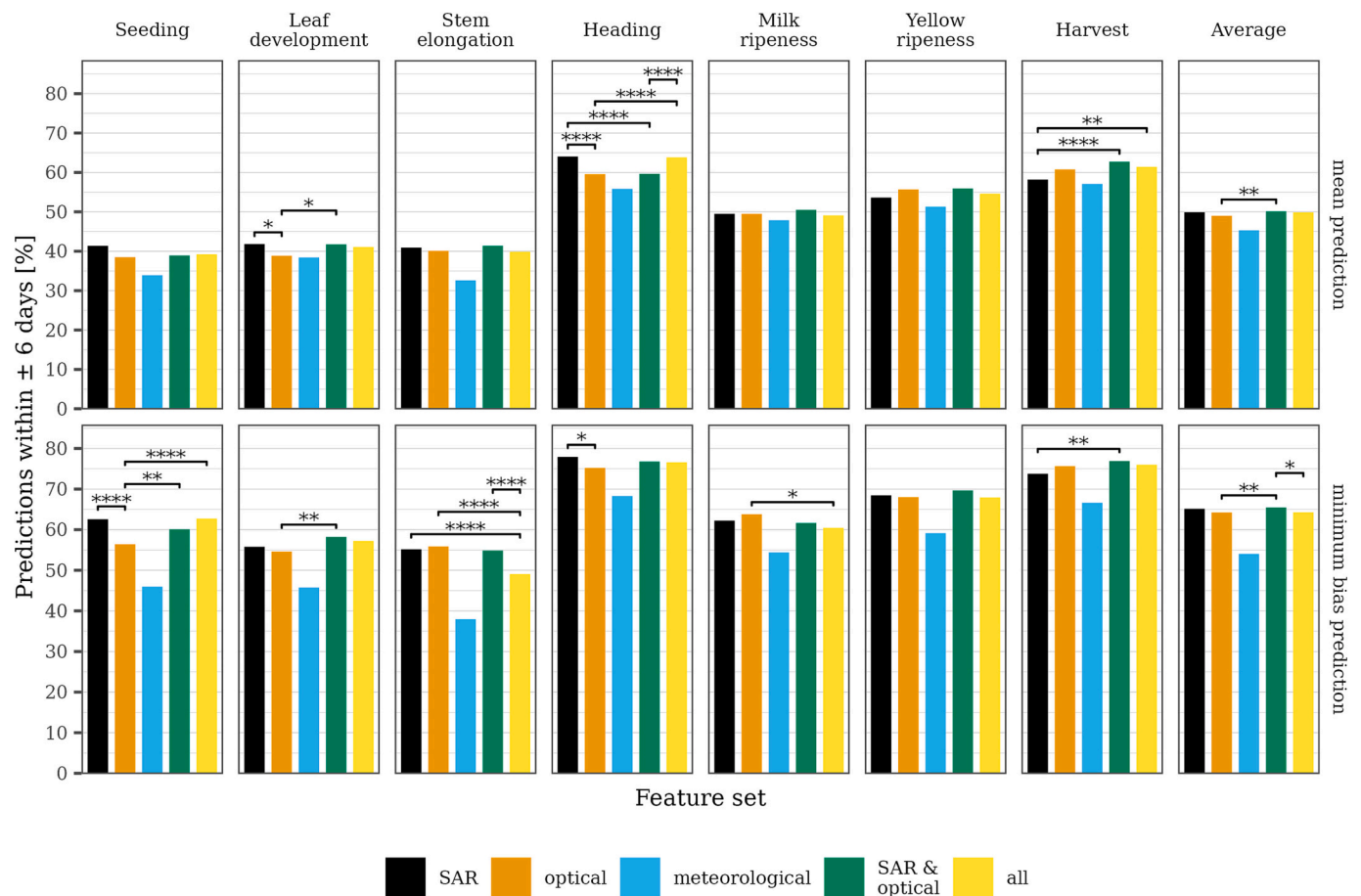


Fig. 11. Prediction accuracy separated by feature set and phenological stage. Brackets indicate significant differences between feature sets according to McNemar's test (McNemar, 1947). All comparisons with the meteorological feature set were significant and therefore excluded for improved readability. Significance was classified as follows: *: $p \leq 0.05$, **: $p \leq 0.01$, ***: $p \leq 0.001$, ****: $p \leq 0.0001$.

optical data as the best feature set and focused our further evaluation on it. An example prediction of the model based on the SAR and optical feature set is included in the appendix (Fig. A1).

4.2. Model baseline

In Fig. 12, we compare the baseline RF model trained on the SAR and optical feature set with our one-dimensional U-Net model. While the RF model showed better *mean* predictions for all phenological stages, only three stages (seeding, heading, and harvest) exhibited significant differences based on McNemar's test. However, except for heading, the one-dimensional U-Net model significantly outperformed the RF model in *minimum bias* predictions across all stages. This significant and consistent advantage in minimum bias predictions led us to choose the one-dimensional U-Net model for further analysis. This decision was further supported by the relevance of minimum bias predictions in our study, as they may better account for the nature and associated uncertainties in the reference data.

4.3. Evaluation of phenological estimates

The predicted start of the phenological stages based on the SAR and optical feature set and the MAE and R^2 regarding the reference data are shown in Fig. 13. Among the seven phenological stages, the predictions for harvest agreed best with the reference data. For both *mean* and *minimum bias* predictions, harvest showed the highest R^2 (0.51 and 0.62) and lowest MAE (5.3 and 4.4). The stage of heading, which reached the highest prediction accuracy (see Fig. 11), showed the second-best MAE (5.4 and 4.5 for *mean* and *minimum bias*) while being in the middle range in terms of R^2 (0.21 and 0.35). Predictions for stem elongation correlated least with an R^2 of 0.06 and 0.28 and an MAE of 9.7 and 7.8, for *mean* and *minimum bias* predictions.

In line with the results for the different feature sets, R^2 and MAE generally improved from *mean* to *minimum bias* predictions. For stem elongation and leaf development, R^2 varied most, with 0.06 compared to

0.28 and 0.09 to 0.44, respectively. Stages with higher R^2 for *mean* predictions improved less, e.g., yellow ripeness showed an R^2 of 0.35 and 0.52 for the *mean* and *minimum bias* predictions.

4.4. Exploration of spatial and temporal patterns

We visualized spatial patterns of our predictions and the reference data for two exemplary phenological stages. For the maps, we decided on the stages with the highest and lowest agreement between the reference data and our model predictions, i.e., harvest (Fig. 14) and stem elongation (Fig. 15). Maps for the other stages including difference maps are shown in the appendix (Fig. A2-Fig. A13).

Reference dates for harvest show a general pattern over the four observed years from an earlier harvest in the south of Germany to a later harvest in the north. Besides this general gradient, we identified regional patterns. An example here is the upper Rhine Valley in the southwest along the border to France, which showed a comparably early harvest in both the reference data and the predictions for the four studied growing seasons.

Our model was able to reproduce these patterns both in the *mean* and *minimum bias* predictions. Trends on a national scale (e.g., overall earlier harvest in 2018) were also reproduced by the model. An evident difference was the significantly higher fine-scale variation in the spatial patterns of the reference data compared to the model predictions. The *minimum bias* predictions better reflected this variation. Yet, both predictions were much smoother and showed remarkably less variance. Predictions for stem elongation showed similarly smooth patterns and a longitudinal gradient, albeit weaker. Yet, reference data showed a much higher level of variance and hardly any trend or pattern for this stage.

We further analyzed the temporal distribution of both the model predictions and the reference data (Fig. 16). The distributions provided insights into the model's capabilities to cover the full temporal spectrum of the reference data. For harvest and heading, e.g., the distribution of the predictions matched the reference data well during nearly all four years. Milk ripeness and yellow ripeness also resembled the

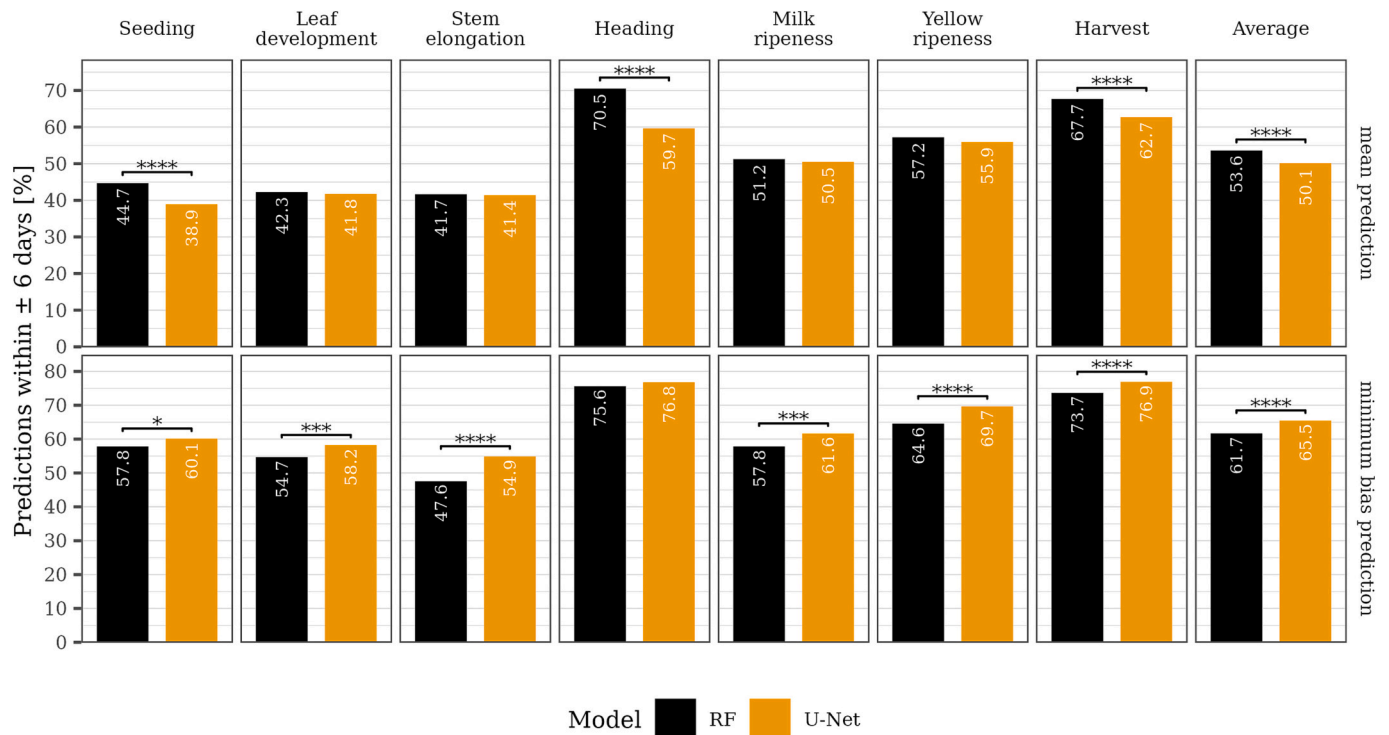


Fig. 12. Prediction accuracy for the RF and U-Net models based on the combination of SAR and optical data. Brackets indicate significant differences between the models according to McNemar's test (McNemar, 1947). Significance was classified as follows: *: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$, ****: $p < 0.0001$.

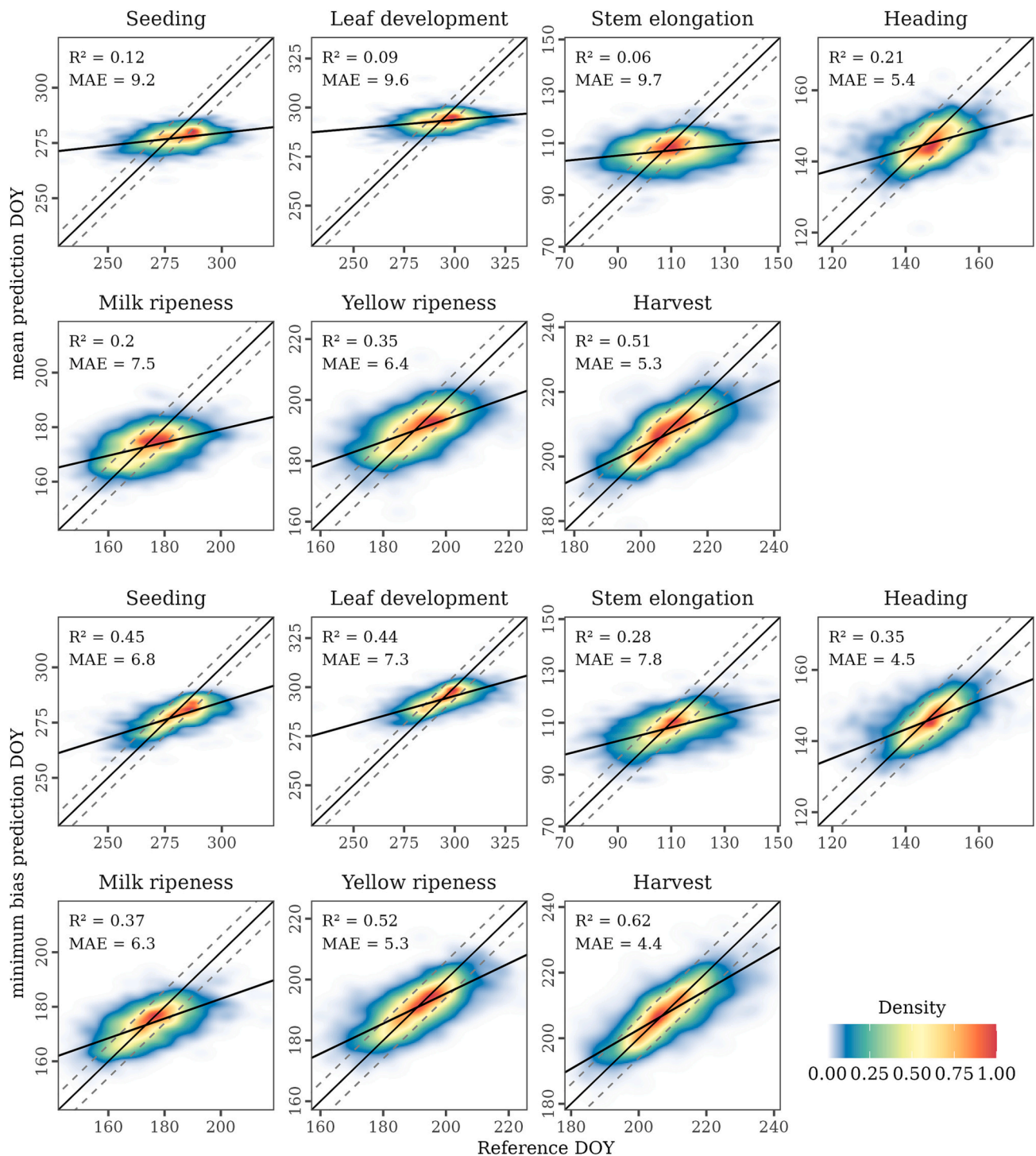


Fig. 13. Density plots of the predicted start of the phenological stages and corresponding reference data for all years based on the combined optical and SAR feature set. Solid lines give the identity (prediction = reference) and regression line. Dashed lines show a deviation of ± 6 days from a perfect prediction, which corresponds to the prediction accuracy reported in Section 4.1.

distributions, but with gaps towards the extremes of the distribution. For seeding, leaf development, and stem elongation the model predictions had less variation and larger gaps.

The leave-one-year-out cross-validation showed differing results for the minimum bias predictions between four analyzed years (Fig. 17). On average, transferring the model to 2017 and 2019 did not show a difference, 2020 showed an overall above-average, and 2018 an overall reduced R^2 . When looking into the individual phases, however, more details can be found. Remarkable is, e.g., the decreased performance for

the phases seeding, heading, and yellow ripeness in 2018 and leaf development and stem elongation in 2019. Next to a decrease, we could also observe higher performance for leaf development and stem elongation in 2020 and seeding and milk ripeness in 2019.

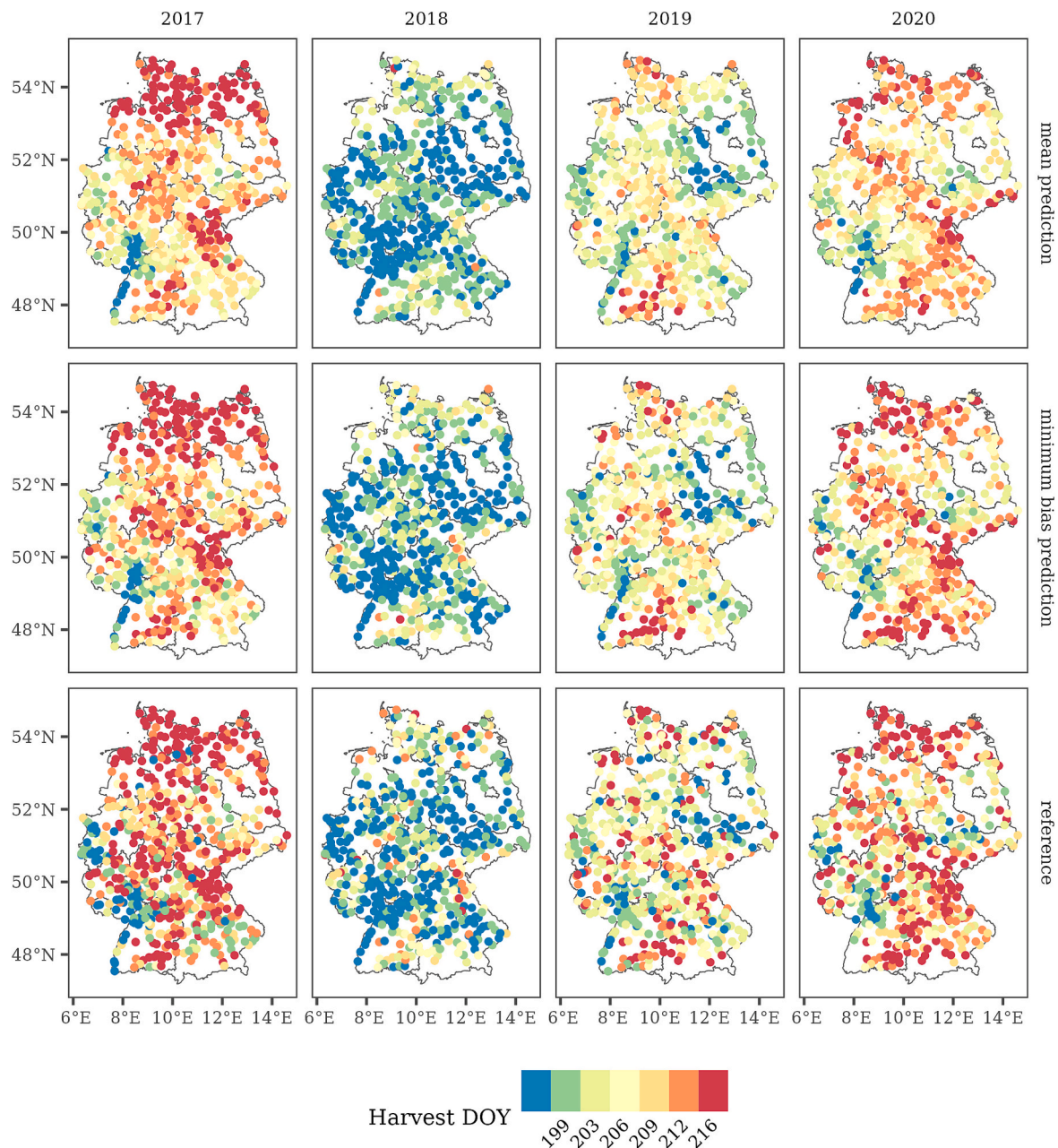


Fig. 14. Maps of predicted and reference dates for the harvest of winter wheat in Germany between 2017 and 2020.

5. Discussion

5.1. Comparison of input features

5.1.1. Individual input features

We evaluated the performance of different remote sensing input features for deriving field-level phenology for winter wheat. Here, we directly aimed to estimate the start of specific phenological stages. This approach distinguishes our study from the common approach of calculating phenological metrics from time series and comparing them with phenological field measurements - sometimes even on a highly aggregated level. Therefore, a direct comparison with other studies is not always straightforward.

The tested feature sets showed different performances. Yet no feature set significantly outperformed the others across all phenological stages. On average, SAR data only performed slightly better than optical data. This supports the findings by [Meroni et al. \(2021\)](#) who compared

different LSP metrics derived from S1 and S2 for winter cereals to aggregated phenological observations from DWD. For winter wheat, they found better agreement with the ground observations for metrics derived from S1 backscatter cross-ratio compared to S2 NDVI. Yet, their overall conclusion was that SAR and optical data perform similarly well, which resembles our findings. [Mercier et al. \(2020\)](#) reported different findings. They used data from both S1 and S2 and compared several optical vegetation indices as well as backscatter coefficient and polarimetric indices to map phenological stages of winter wheat targeting eight acquisition dates. In opposite to our results, they found optical data to yield higher accuracies compared to SAR data. However, their approach considerably differs from our work since it completely omits the temporal domain of satellite data. Other studies reported differences between the performance of optical and SAR data, yet did not reach a general conclusion (e.g., [Harfenmeister et al., 2021](#); [Velooso et al., 2017](#)).

Looking at the individual stages, radar data tended to work better for earlier stages (seeding to heading). This is consistent with the

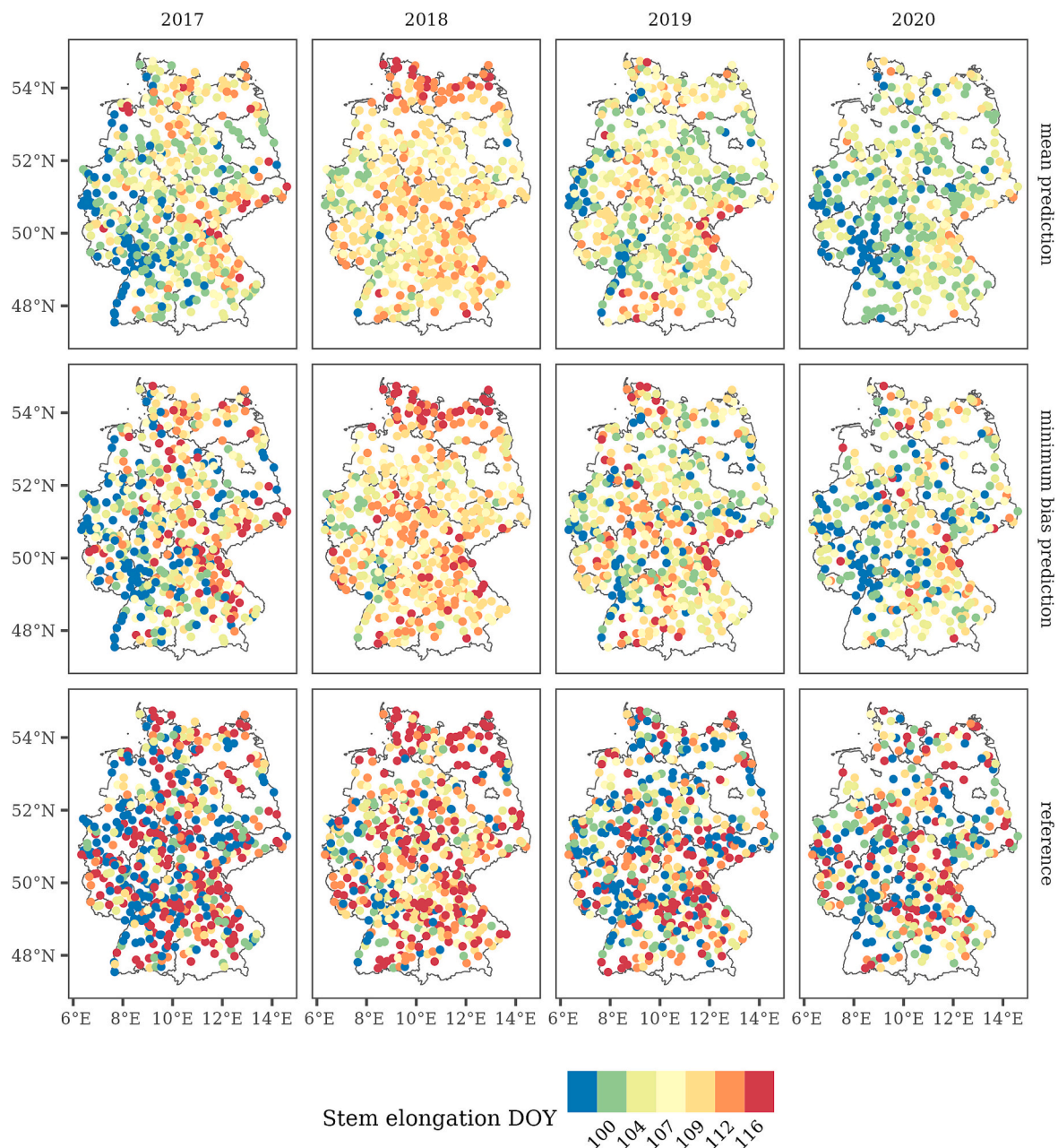


Fig. 15. Maps of predicted and reference dates for the stem elongation stage in Germany between 2017 and 2020.

observations of Jia et al. (2013) who conducted a ground-based radar backscattering experiment in different frequencies and polarizations for different phenological stages of winter wheat. Overall, they found the backscatter coefficient to be more sensitive to changes during the early growing period, followed by a decline towards the maturity of the crop. For seeding, the superior performance of SAR data might be due to the sensitivity of the SAR signal for soil roughness. Seeding is usually closely accompanied by tillage practices, that significantly change the soil structure and hence the SAR signal, while the multispectral, optical signal might experience less change. Another potential reason for SAR being more sensitive to subtle changes might relate to the high observation density compared to optical data during this time. Seeding is usually done in fall, which is a season with frequent cloud cover in Germany (see Fig. 6).

SAR data also outperformed optical data for leaf development. This supports previous findings on C-band SAR data for detecting thin wheat seedlings, even though different incidence angles will yield different

results (Jia et al., 2013). Fieuzal et al. (2013) reported SAR being more sensitive to stem elongation compared to NDVI. Here, our results are not clear and show differences between *mean* and *minimum bias* predictions. For heading, differences in SAR and optical data can mainly be explained by structural changes of the wheat plant during this period. The heads emerging from the leaf sheet may have less influence on the spectral signature compared to the SAR backscatter. Meroni et al. (2021) also raised this hypothesis after observing a clearer signal in the backscatter cross-ratio compared to the NDVI during heading.

For later phenological stages (milk ripeness to harvest) we conversely found that optical data outperformed radar data. While this was also stated by Mercier et al. (2020), who observed optical data being better suited for detecting the end of ripening, Meroni et al. (2021) found the opposite for the stage of yellow ripeness. The transition from milk ripeness to yellow ripeness comes with clear changes in color as well as a decline in photosynthetic activity. Both highly influence the spectral signature and could explain the advantage of optical data.

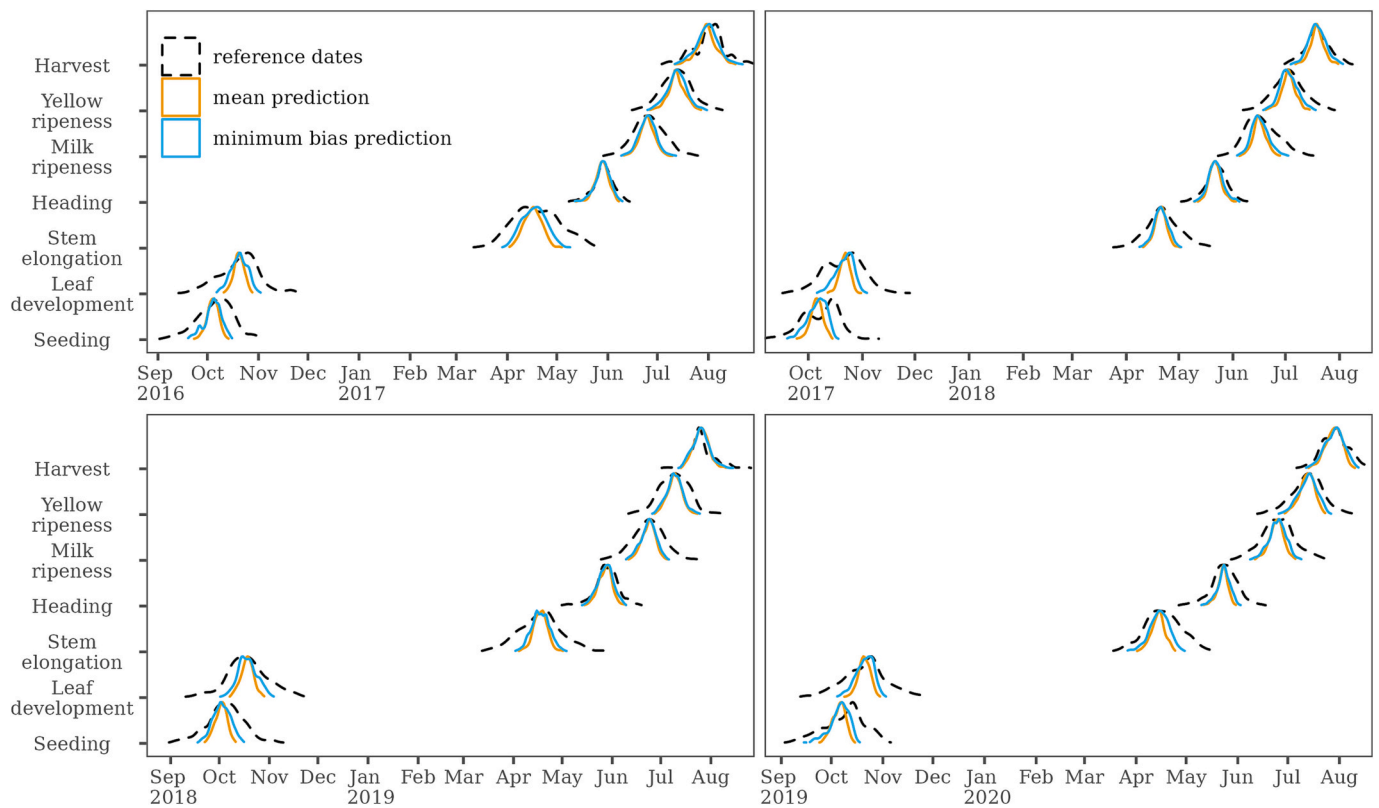


Fig. 16. Temporal distribution of the predicted start of the phenological stages and corresponding reference data.

However, simultaneously the water content decreases during this time which influences the plant's dielectricity and hence may also influence the SAR signal.

We did not observe a significant difference between optical and SAR data for detecting the harvest, although optical data showed better accuracies on average. Meroni et al. (2021) observed clearer changes in NDVI compared to cross-ratio around harvest and mentioned the similarity of fully mature plants and stubble in the SAR signal as possible explanation. This may highly depend on the harvesting practice. Other studies, however, do not report such findings (e.g., Fieuzal et al., 2013; Nasrallah et al., 2019). Harfenmeister et al. (2021) argue that the very low photosynthetic activity directly before harvest also does not have to result in major differences in NDVI before and after harvest. Beyond NDVI, we also know that especially Shortwave Infrared (SWIR) reflectance is suitable for distinguishing dry biomass and soil which supports our findings (Daughtry, 2001). Using coherence data could further improve harvest detection, as it was observed to be useful to detect the harvest of cereals and mowing of grasslands in other studies (Kavats et al., 2019; Lobert et al., 2021).

The meteorological variables showed the least explanatory power among our tested input features. Gerstmann et al. (2016) demonstrated that an approach based solely on meteorological data yielded great explanatory potential for the timing of crop phenological development. Yet, they studied phenological development on a 1 km² grid size. In our study, we could identify high variation of phenological development on a fine scale from the mapped reference dates for harvest (Fig. 14) and stem elongation (Fig. 15). Micro-relief, soil properties, and management practices are potential influencing factors. Meteorological data, especially of the spatial resolution used in our study, cannot represent these variations. For example, the timing of management may vary vastly between fields belonging to an in-situ observation, while the meteorological conditions can be similar. This also becomes evident from the small increase in performance when comparing *mean* and *minimum bias* predictions for the meteorological feature set, which indicates that the

individual predictions for the fields belonging to the same observation create similar predictions.

5.1.2. Feature combinations

Combining SAR and optical data did not significantly improve the model performance in our study. Even if we observed an increase in prediction accuracy over one of both sensors alone, we could not clearly confirm the findings by Mercier et al. (2020), who found an improvement by combining S1 and S2 data for their phenology classification algorithm. This synergy was also suggested by Harfenmeister et al. (2021), Veloso et al. (2017), and Yeasin et al. (2022) who found vastly different but also complementary performances of SAR and optical time series for analyzing the phenology of winter wheat, barley and sugarcane. The improvement for some stages could be attributed to uncertainties and ambiguities in the predictions with SAR or optical data alone, respectively, that could be resolved by combining both. When precipitation-induced noise in the SAR signal or data gaps in the optical data hamper the precise delineation of the event in the time series, combining both enables our proposed model to refine the predictions.

For some stages, a decrease in accuracy was observed when combining optical and SAR data. Including data with low or redundant information content can make it harder for the model to identify patterns in the increased amount of data. The noise introduced by such data can hinder the model's ability to extract meaningful information, leading to decreased performance and accuracy (Bellman, 2003). This hypothesis is supported by our observation that the largest decrease occurred when the performance difference between single-sensor (optical, radar) feature sets was particularly large, i.e., *mean* predictions for heading and *minimum bias* for seeding.

Considering the combination of remote sensing imagery with meteorological information, we observed an increase in model performance for the seeding stage. When heavy rainfall events have just occurred or are forecasted, the farmer might reschedule the seeding date due to, e.g., non-accessible soils. Including such information could have enabled the

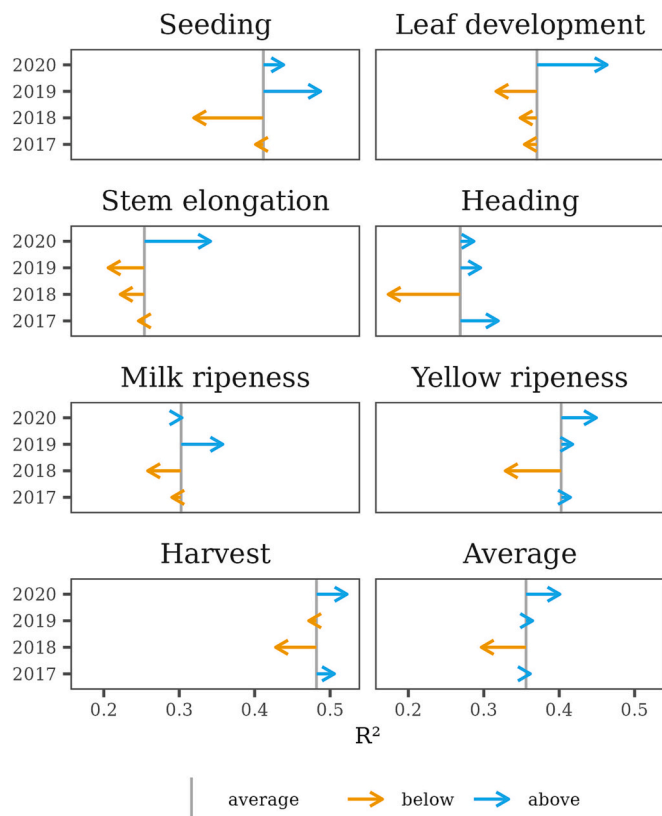


Fig. 17. Results of the leave-one-year-out cross-validation for the seven phenological stages as well as the average for all stages for minimum bias predictions. Vertical line shows the mean for the respective phase over the four years. Directions of the arrows are indicating if the performance for a specific year was above or below the average of all four years, and lengths are indicating the magnitude of the difference.

model to account for such events. For the heading stage, adding meteorological to SAR and optical data also improved the predictions. This matches the findings from Gerstmann et al. (2016) who observed the best performance for the heading stage compared to other stages using meteorological data only. However, only the *mean* predictions improved for heading, which does not support a high explanatory power for field-based estimates, since the resolution of the temperature data we used provides only limited variations between the fields. Overall, our results suggest that meteorological data do not add significant value to dense remote sensing time series for phenological monitoring.

5.2. Model baseline

While the one-dimensional U-Net model demonstrated significantly superior performance for the *minimum bias* predictions, we have also seen Random Forest to achieve similar to even better results in *mean* predictions. Although we put more weight on the minimum bias predictions for the model choice and thus decided for the U-Net, this nevertheless demonstrates the potential that already exists in state-of-the-art machine learning algorithms for phenology analysis. However, besides accuracies and statistical significance, it is essential to also consider the practical implications of model choice. Our chosen U-Net architecture provides a significant advantage for phenology monitoring by offering highly detailed output with assigned probabilities for each time step, indicating the start of the seven phases (see Fig. 10). This multidimensional granularity enables comprehensive research into winter wheat's phenological development. In comparison, models like Random Forest typically predict a single target value per input sample. Replicating the U-Net's output using alternative models would require

training multiple models and implementing auxiliary preprocessing steps, such as generating moving windows. This approach would be time-consuming and prone to errors. In contrast, the U-Net architecture offers an efficient and streamlined solution for full-season phenology predictions without the need for an extensive ensemble of models or complex preprocessing, which is especially important for long-term monitoring tasks.

5.3. Evaluation of phenological estimates

For the best model (SAR & optical), prediction accuracies increased from early towards later phenological stages. This is in line with Gerstmann et al. (2016). Later stages are associated with almost complete plant coverage. They are accompanied by significant structural changes (e.g., heading), vast changes in color and water content (yellow ripeness), or a combination of changes (harvest). These changes affect signals from both optical and radar sensors and indicate good detectability. Milk ripeness shows less obvious or abrupt changes that could be detectable by SAR or optical sensors, which is also reflected by the relatively low R^2 compared to the other late stages. Zeng et al. (2020), however, reported that the estimation of phenology information during the vegetation's senescence is a greater challenge compared to the green-up. Comparable limitations for later stages were also reported by Harfenmeister et al. (2021) and Shang et al. (2020) for SAR-based methods.

During the early stages crop cover is not present at all (e.g., seeding) or is still low (e.g. stem elongation). This leads to a high proportion of soil signal in the remote sensing imagery and only little signal attributable to vegetation. Despite the aforementioned sensitivity of radar data to small seedlings or tillage, these stages apparently provide less distinctive features in the time series that could be recognized by our model.

For all feature sets and phenological stages, we have seen an increase in the prediction accuracy, a decrease in MAE, as well as an increase in R^2 from the *mean* towards the *minimum bias* method. This increase was also observed by Ye et al. (2022). Especially for the stages with the highest differences (e.g., seeding), this observation indicates that our model predictions cover some temporal range - even between the candidates for one field observation - and can also predict the phenology of fields that differ from the mean in a given area.

The low slope of the regression lines indicates that the predictions do not cover the full temporal range and variation of the reference data. However, the aggregated nature of the predictions also plays a role here, making them more likely to tend towards the average of a region as opposed to the reference data, which comes from only one field and thus could be both representative or an outlier compared to the surrounding fields.

5.4. Exploration of spatial and temporal patterns

Mapping the predictions and reference dates for the phenological stages provided us with valuable insights into their spatial distribution. The consistent pattern of the predictions throughout the years indicates that our method generates regionalized results that reflect the overall environmental conditions in Germany well.

The similarity between the distributions of predicted and reference dates for later phenological stages indicates that our model covers both the spatial and temporal gradients of these stages across Germany. The model could, therefore, also predict fields where the phases began sooner or belated. This finding also suggests that the model is well suited for the area-wide prediction of these stages in Germany. For earlier phenological stages, the model's limitation in covering the temporal distribution of the reference data may indicate that the model predicts also based on seasonal trends. This explains the concentration of the distribution towards the distribution means (Fig. 16) and the narrow ranges of estimation (Fig. 13).

For stem elongation, the high level of variance in the reference data

could not be reproduced. This may be explained by the combination of different sensor types still not providing sufficient information to precisely detect such subtle variations in spectral or backscatter behavior. Another factor for the limited predictions could be the reference data. On the one hand, these could be affected by uncertainties (compare Section 5.5). On the other hand, the sampled field itself could be a statistical outlier compared to the surrounding fields, which is difficult to account for with our methodology.

The leave-one-year-out cross-validation revealed the temporal transferability of our proposed model in dependence on the individual phenological stages. The decreased below-average performance when leaving out 2018 could be explained by exceptional weather conditions (see Fig. 2). Starting with wet conditions during seeding and leave development in 2017, 2018 started with a relatively cold period followed by comparably high temperatures and little precipitation for the whole vegetation period. While this reduces the impact on optical time series through scarce cloudiness and SAR time series through low soil moisture influence, the whole phenological timing was exceptional in that year as becomes apparent from Fig. 3. In contrast, shorter dry periods as in June and July 2019 show above-average performance exemplified by milk ripeness that occurred at that time. The same applies to seeding in 2019. However, above-average performance for, e.g., leaf development in 2020 cannot solely be explained by weather phenomena and suggests that other factors are also influencing the model predictions.

5.5. Limitations and outlook

We based our study on a reference data set from a national phenology network. As discussed in detail by Ye et al. (2022), such data have their strengths but also provide some challenges. While covering broad geographical and ecological extents, using such data for training and validating predictive models might be hampered by noise and errors in the reference data related to the way observers report phenology. Such a volunteer-based approach may result in differences between the actual and reported start of the stages if volunteers are not visiting fields on a daily basis. For example, the “weekend bias” is a known phenomenon described by Courter et al. (2013). Furthermore, even if the observers are trained, misclassifications of phenological stages are possible.

A major challenge discussed by Ye et al. (2022) is the missing link between in-situ observations conducted on a single plant or field and mixed pixels in remote sensing data. Using an LSP product with 500 m spatial resolution from the Visible Infrared Imaging Radiometer Suite (VIIRS), Ye et al. (2022) suggested several methods to upscale multiple in-situ observations to the VIIRS pixels. We adopted and inverted this approach to aggregate multiple field predictions to match one in-situ observation using the *mean* and *minimum bias* methods. Using both methods, we were able to validate our model predictions and also gain insights into prediction variations by comparing the results of both methods. Our approach was well-suited as a reference for comparing different input features. However, metrics resulting from the *minimum bias* predictions should be interpreted with caution, as they are not completely independent from the reference data (Ye et al., 2022).

Our field boundary generation allowed us to relate the field measurements to field-based remote sensing time series. However, two differently managed, neighboring winter wheat fields may be lumped into the same boundary. Especially for management-related stages, i.e., seeding and harvest, this can lead to a mixture of temporal profiles, where patterns for the corresponding stages could occur twice or blend into each other. A possible solution for future work would be the use of more sophisticated field delineation approaches that can account for management practices (e.g., Tetteh et al., 2021).

The proposed method using DL enabled us to combine and simultaneously exploit time series of different remote sensing sensors and meteorological measurements. The great flexibility of DL models enables to adapt their architecture to any given problem. Here, it allowed

us to predict the start of several phenological stages at the same time based on a variety of feature sets and assess the performance of different combinations with a single streamlined model. Further research should focus on extending model architectures with a spatial dimension and testing more data sources that provide additional information (e.g., coherence) or come with higher spectral or temporal resolution.

6. Conclusion

We demonstrated the overall capability of a one-dimensional temporal U-Net model to simultaneously predict the start of the major phenological stages for winter wheat based on SAR and optical remote sensing time series for individual fields. Even if we observed an increase in accuracy our results could not undoubtedly confirm the synergistic potential of optical and SAR remote sensing data for such purposes. We also did not find a general improvement in our results when adding meteorological variables to the model. We conclude that precipitation data (e.g. from a rainfall radar network) or interpolated temperature measurements alone are not able to explain fine-scale differences of phenology at the field level that are rather related to farmers' decisions on cropping practices. The strengths of radar data especially supported analyses at the earlier stages of plant development between seeding and heading. After the complete formation of the stand and in the subsequent phases of maturity and senescence, the optical data gained importance.

This study is a step forward towards directly targeting explicit phenological stages when dealing with vegetation analyses from remote sensing data. Despite well-known limitations of national-scale phenological observations, we proposed a calibration scheme that enables to combine such data with field level time series. Although we were able to make field-level predictions, a field-level validation in the strict sense was not possible here. However, based on an adapted validation strategy, we were able to valorize the unique and Germany-wide phenology reference dataset and to underline the additional value and necessity of field-level reference data for future model optimization. We further demonstrated that DL models provide great flexibility that allows adapting them to a broad range of problems and tasks.

Overall, this study adds to our knowledge base on remote sensing-based high-resolution mapping of vegetation productivity from space. The proposed method is ready to be applied for area-wide assessments of vegetation phenology at the national level and beyond. It can next be tested for investigating management-related influences on crop phenology at the field level and, thus, cropland use intensity. Ultimately, it may be used for the evaluation of agricultural and environmental policies.

CRediT authorship contribution statement

Felix Lobert: Conceptualization, Methodology, Software, Formal analysis, Investigation, Data curation, Visualization, Writing – original draft. **Johannes Löw:** Investigation, Writing – review & editing. **Marcel Schwieder:** Software, Methodology, Investigation, Writing – review & editing. **Alexander Gocht:** Writing – review & editing. **Michael Schlund:** Writing – review & editing. **Patrick Hostert:** Conceptualization, Supervision, Writing – review & editing. **Stefan Erasm:** Conceptualization, Methodology, Supervision, 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.

Data availability

Data will be made available on request.

Appendix A

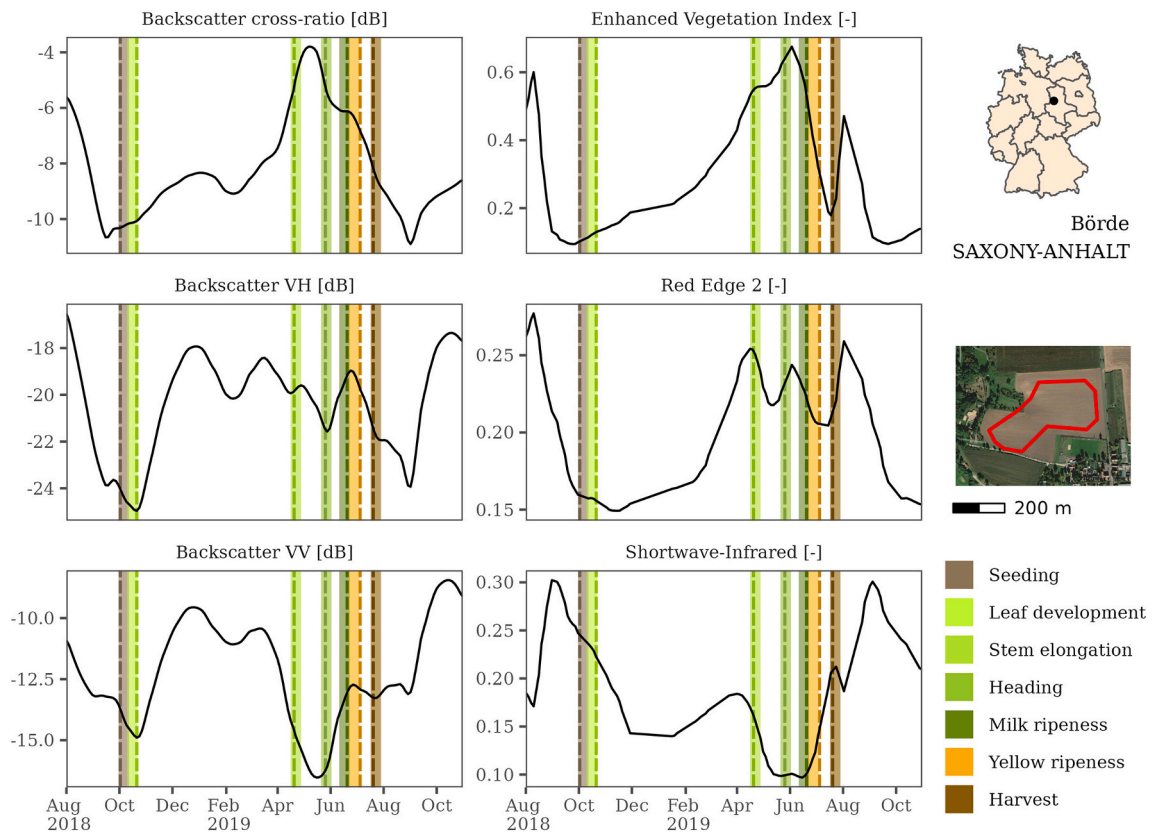


Fig. A1. Predicted start of the different phenological stages for an exemplary winter wheat field with a selection of optical and SAR-based input features. Dashed vertical lines show the prediction, segments in the background give the reference date including a buffer of 6 days.

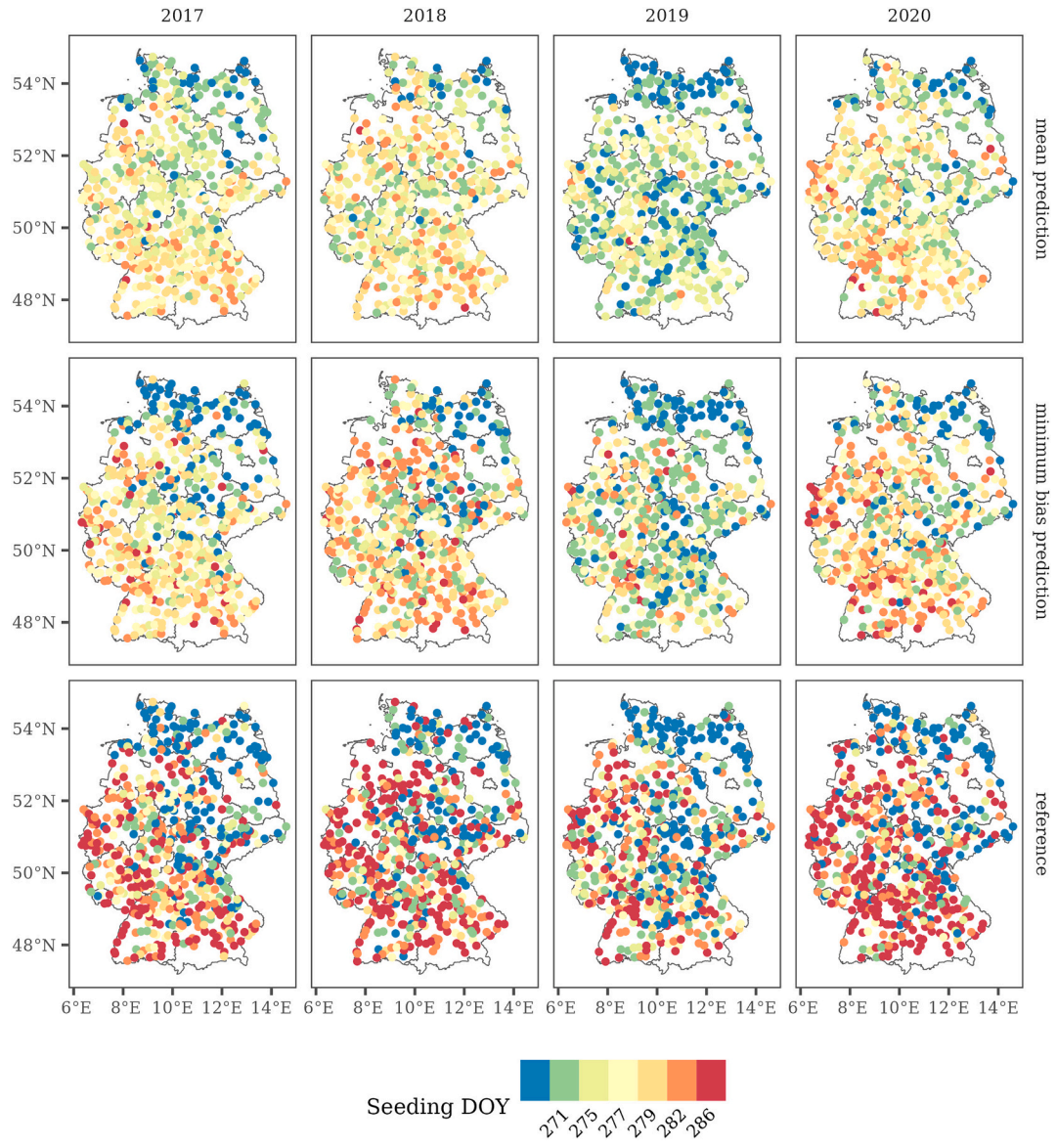


Fig. A2. Maps of predicted and reference dates for the seeding of winter wheat in Germany between 2017 and 2020.

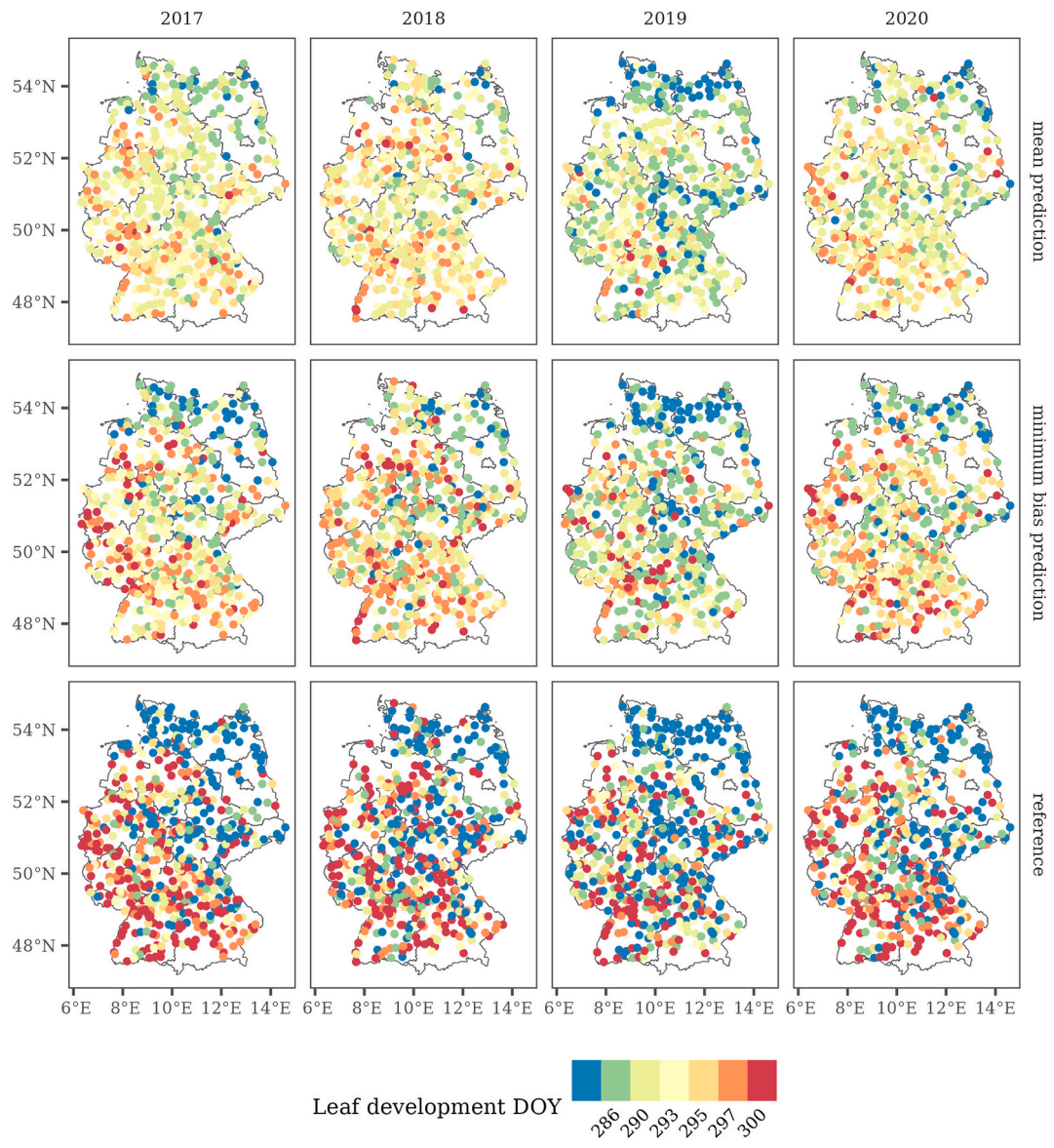


Fig. A3. Maps of predicted and reference dates for the leaf development of winter wheat in Germany between 2017 and 2020.

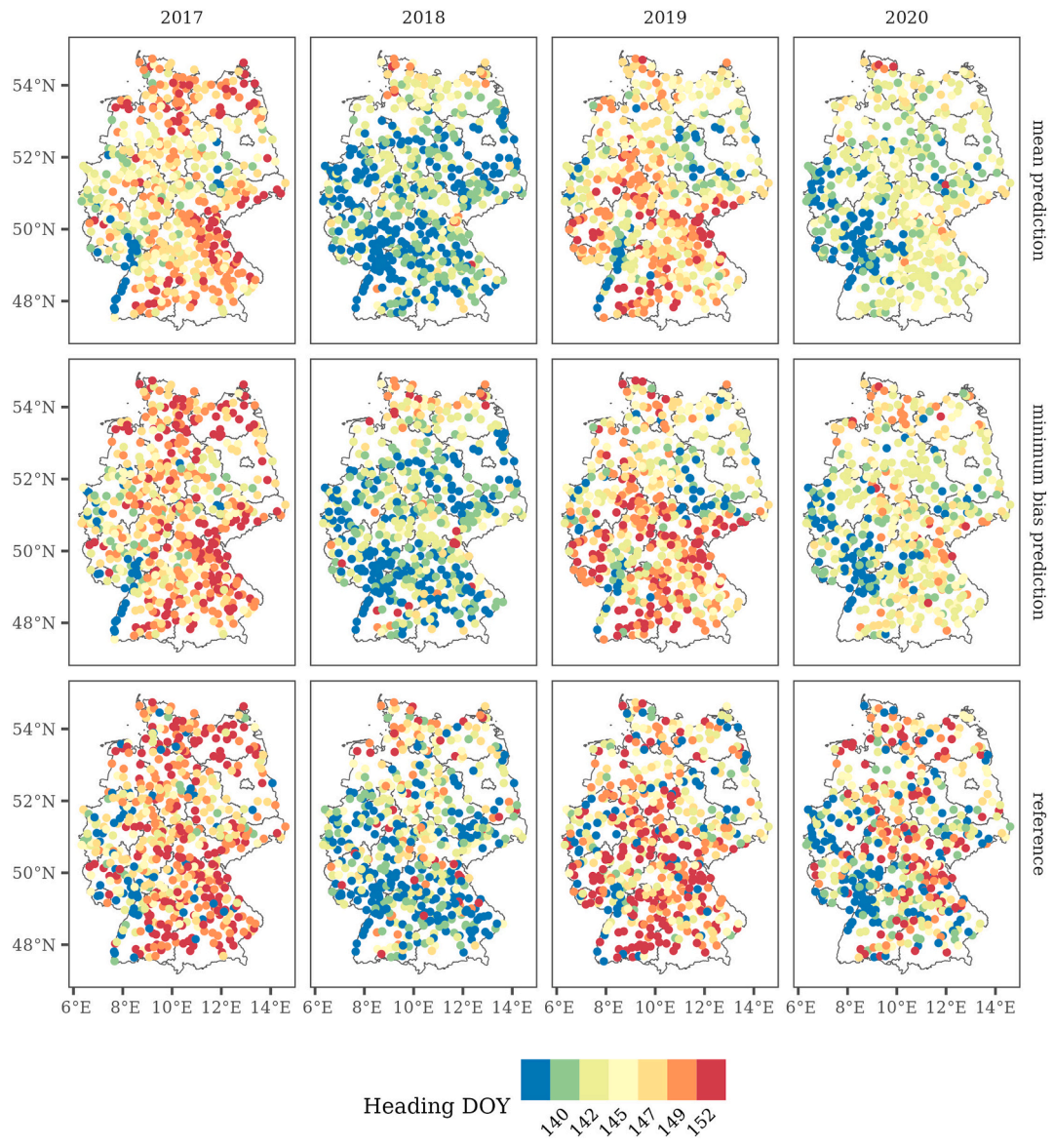


Fig. A4. Maps of predicted and reference dates for the heading of winter wheat in Germany between 2017 and 2020.

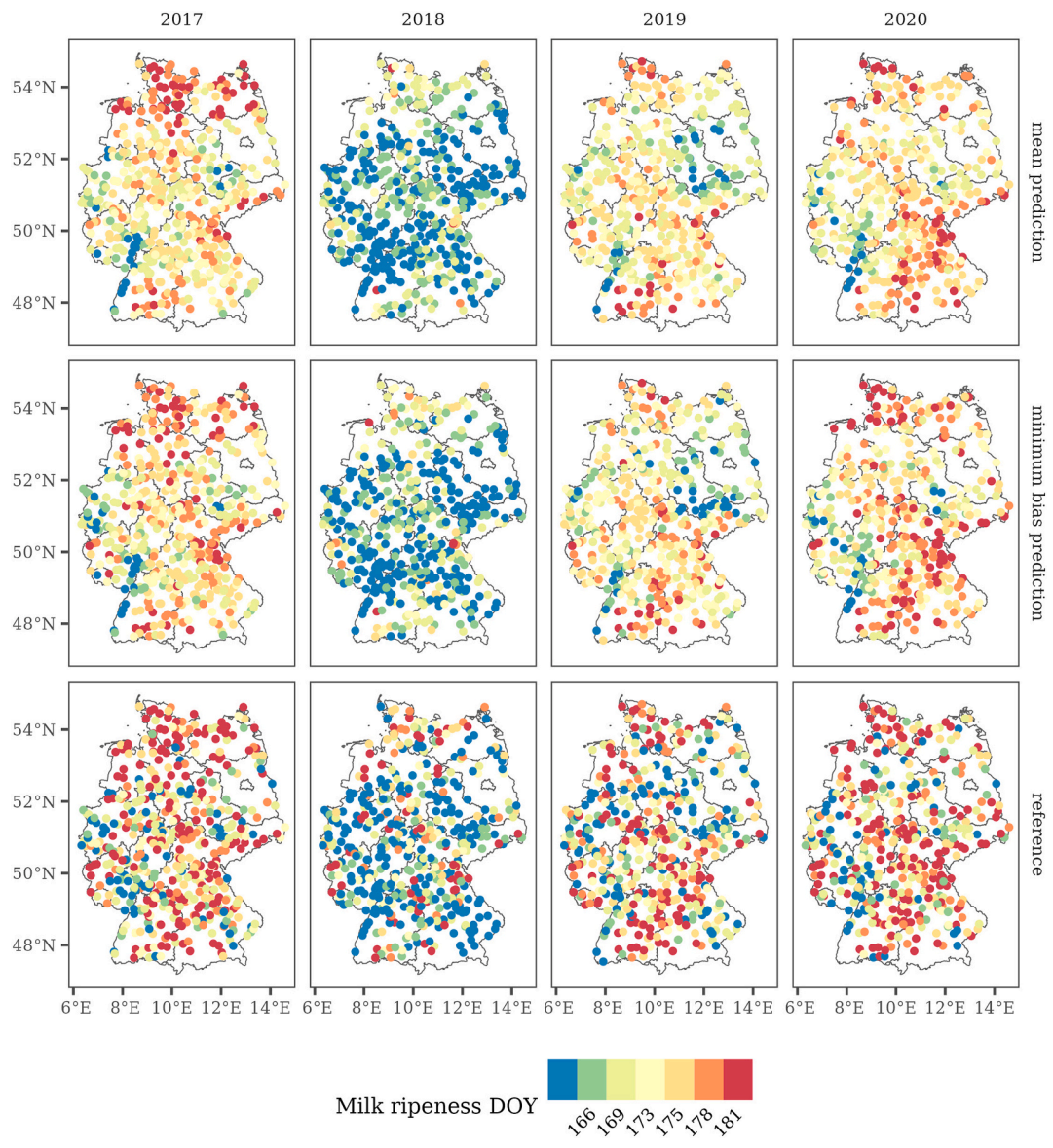


Fig. A5. Maps of predicted and reference dates for the milk ripeness of winter wheat in Germany between 2017 and 2020.

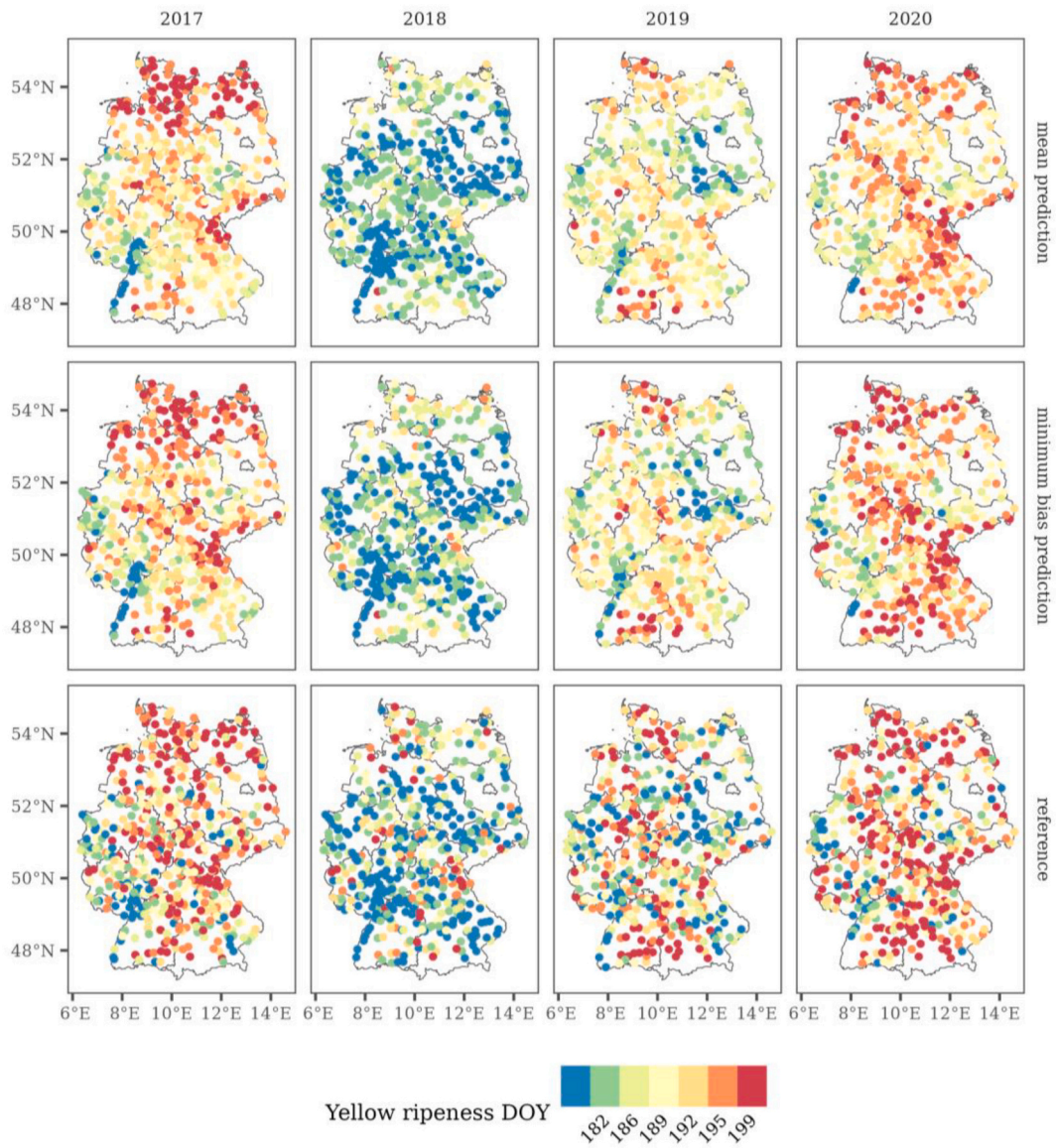


Fig. A6. Maps of predicted and reference dates for the yellow ripeness of winter wheat in Germany between 2017 and 2020. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

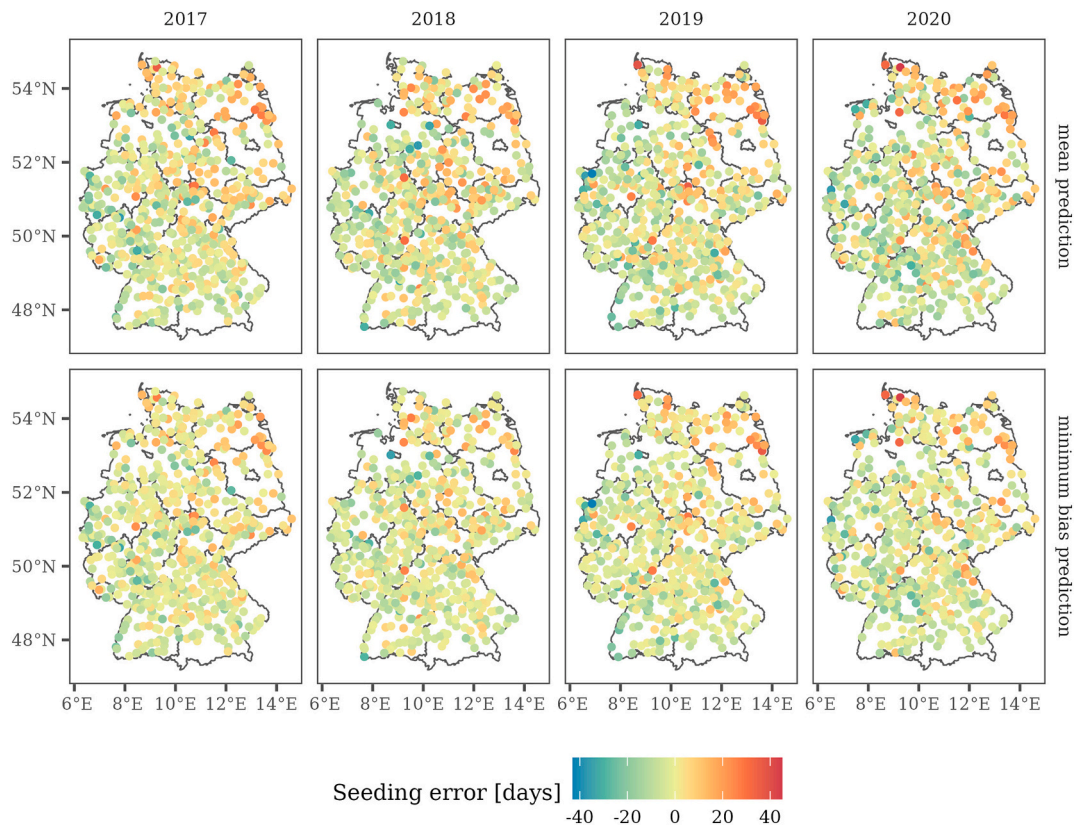


Fig. A7. Difference of predictions and reference dates for the seeding of winter wheat in Germany between 2017 and 2020.

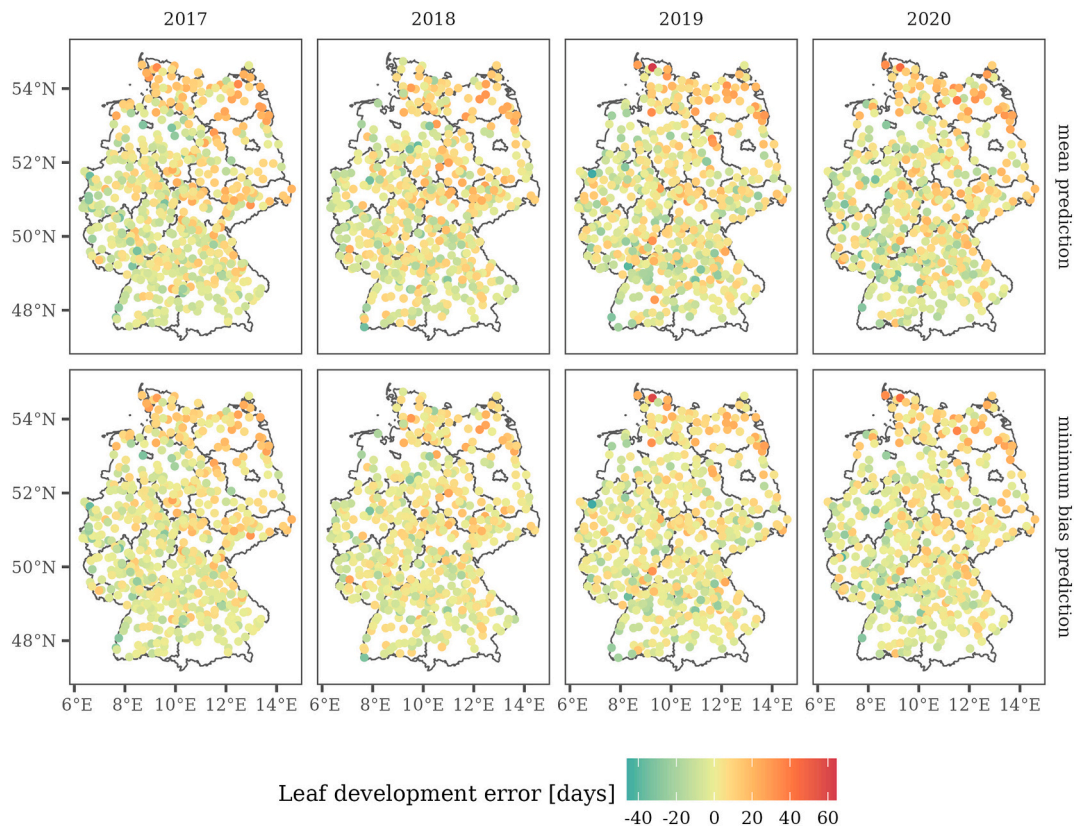


Fig. A8. Difference of predictions and reference dates for the leaf development of winter wheat in Germany between 2017 and 2020.



Fig. A9. Difference of predictions and reference dates for the stem elongation of winter wheat in Germany between 2017 and 2020.

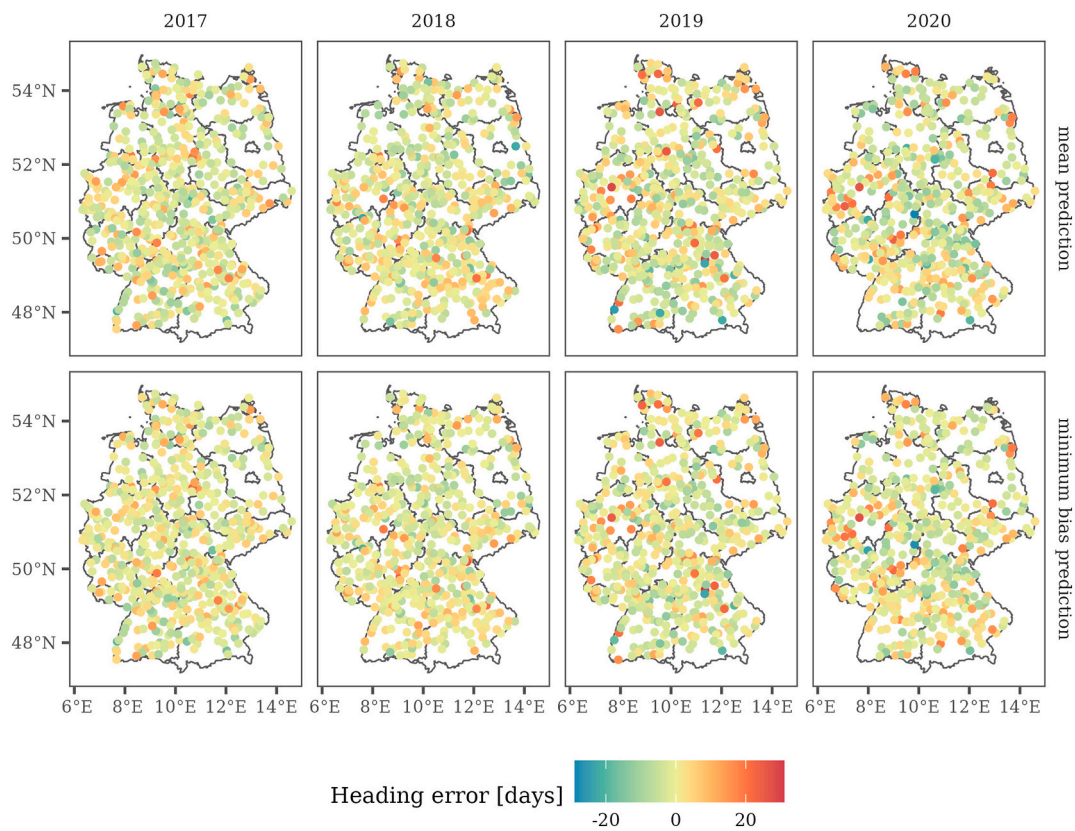


Fig. A10. Difference of predictions and reference dates for the heading of winter wheat in Germany between 2017 and 2020.



Fig. A11. Difference of predictions and reference dates for the milk ripeness of winter wheat in Germany between 2017 and 2020.

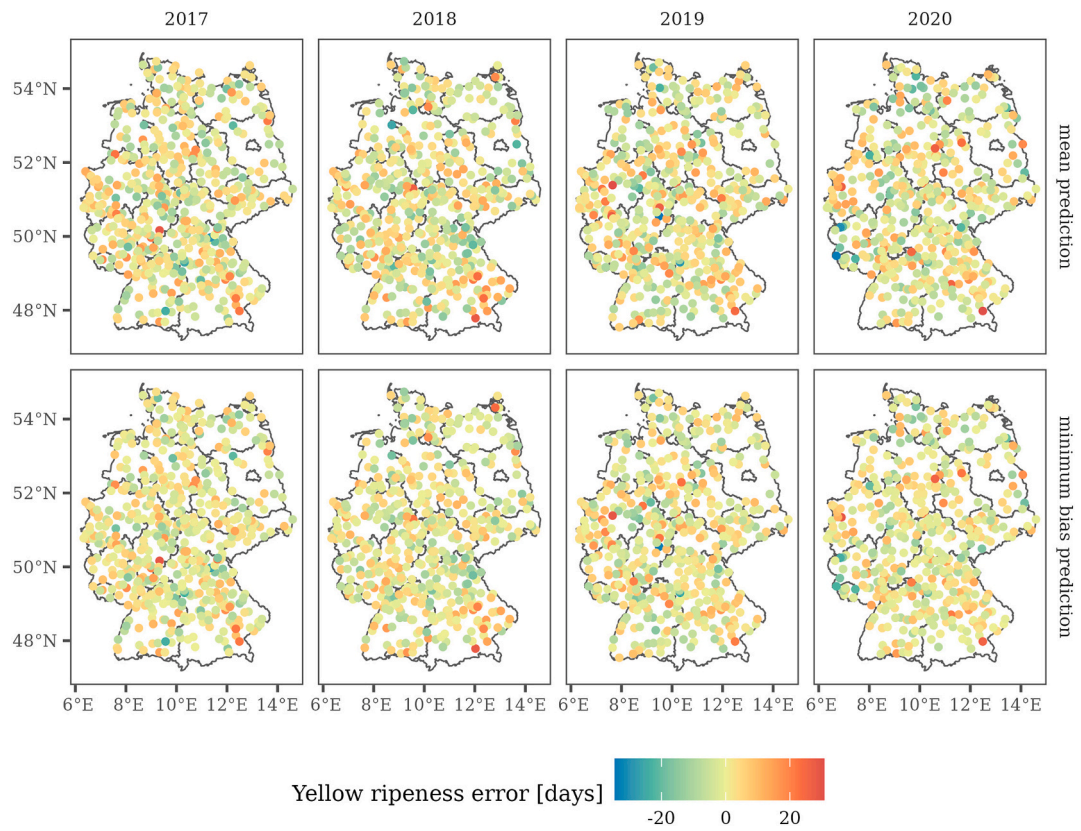


Fig. A12. Difference of predictions and reference dates for the yellow ripeness of winter wheat in Germany between 2017 and 2020. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

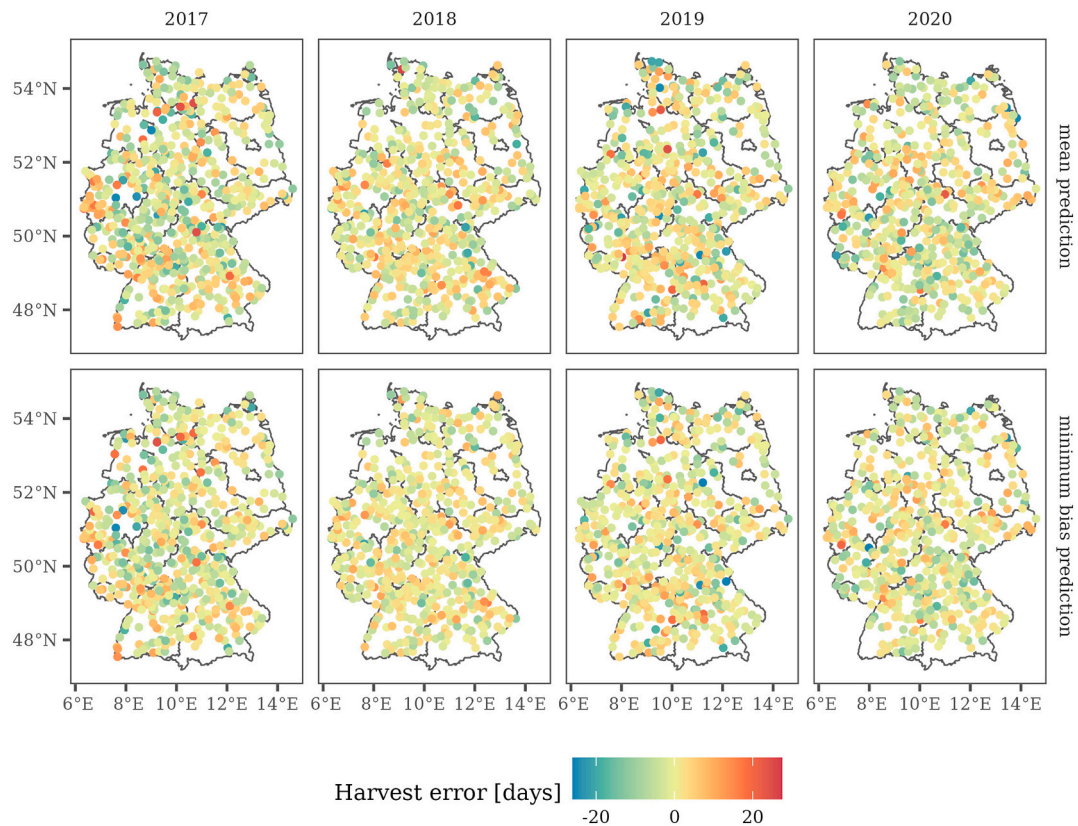


Fig. A13. Difference of predictions and reference dates for the harvest of winter wheat in Germany between 2017 and 2020.

References

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., 2016. In: *Tensorflow: A system for large-scale machine learning*, in: 12th USENIX Symposium on Operating Systems Design and Implementation, pp. 265–283.
- Allaire, J., Chollet, F., 2021. Keras: R Interface to “Keras”. R package version 2 (6), 1. <https://cran.r-project.org/package=keras>.
- Bellman, R., 2003. *Dynamic programming*, 1. publ., repr. of the 6. print. 1972. ed, Dover Books on Mathematics. Dover Publ, Mineola, NY.
- Benz, U., Banovsky, I., Cesarz, A., Schmidt, M., 2020. CODE-DE Portal Handbook, Version 2.0, DLR.
- Bishop, C.M., 1995. *Neural Networks for Pattern Recognition*. Oxford University Press Inc, USA.
- Blickensdorfer, L., Schwieder, M., Pflugmacher, D., Nendel, C., Erasmí, S., Hostert, P., 2022. Mapping of crop types and crop sequences with combined time series of Sentinel-1, Sentinel-2 and landsat 8 data for Germany. *Remote Sens. Environ.* 269 <https://doi.org/10.1016/j.rse.2021.112831>.
- Boessenkool, B., 2021. Rdw: select and download climate data from “DWD”. (German Weather Service).
- Bolton, D.K., Gray, J.M., Melaas, E.K., Moon, M., Eklundh, L., Friedl, M.A., 2020. Continental-scale land surface phenology from harmonized landsat 8 and Sentinel-2 imagery. *Remote Sens. Environ.* 240, 111685 <https://doi.org/10.1016/j.rse.2020.111685>.
- Breiman, L., 2001. Random forests. *Mach. Learn.* 45, 5–32. <https://doi.org/10.1023/A:1010933404324>.
- Chollet, F., 2015. Keras. <https://keras.io>.
- Cleveland, W.S., Grosse, E., Shyu, W.M., 1992. Local regression models. In: *Statistical Models in S*. Routledge.
- Courter, J.R., Johnson, R.J., Stuyck, C.M., Lang, B.A., Kaiser, E.W., 2013. Weekend bias in citizen science data reporting: implications for phenology studies. *Int. J. Biometeorol.* 57, 715–720. <https://doi.org/10.1007/s00484-012-0598-7>.
- d’Andrimont, R., Taymans, M., Lemoine, G., Ceglár, A., Yordanov, M., van der Velde, M., 2020. Detecting flowering phenology in oil seed rape parcels with Sentinel-1 and -2 time series. *Remote Sens. Environ.* 239, 111660 <https://doi.org/10.1016/j.rse.2020.111660>.
- Daughtry, C.S.T., 2001. Discriminating crop residues from soil by shortwave infrared reflectance. *Agron. J.* 93, 125–131. <https://doi.org/10.2134/agronj2001.931125x>.
- De Beurs, K.M., Henebry, G.M., 2004. Land surface phenology, climatic variation, and institutional change: analyzing agricultural land cover change in Kazakhstan. *Remote Sens. Environ.* 89, 497–509. <https://doi.org/10.1016/j.rse.2003.11.006>.
- DWD, 2022. Climate Data Center (CDC). Phenological observations of crops from sowing to harvest (annual reporters, historical), Version v008 https://opendata.dwd.de/climate_environment/CDC/observations_germany/phenology/annual_reporters/crops/historical/PH_Jahresmelder_Landwirtschaft_Kulturpflanze_Winterweizen_1925_2021_hist.tx.
- DWD, 2022. Climate Data Center (CDC). Historical hourly station observations of 2m air temperature and humidity for Germany, version v006, 2018.
- DWD, 2022. Climate Data Center (CDC). Historical hourly sliding RADOLAN grid of daily precipitation (binary), version 2.5.
- DWD, 2015. *Vorschriften und Betriebsunterlagen für die phänologischen Beobachter des Deutschen Wetterdienstes*.
- Farr, T.G., Rosen, P.A., Caro, E., Crippen, R., Duren, R., Hensley, S., Kobrick, M., Paller, M., Rodriguez, E., Roth, L., Seal, D., Shaffer, S., Shimada, J., Umland, J., Werner, M., Oskin, M., Burbank, D., Alsdorf, D.E., 2007. The shuttle radar topography mission. *Rev. Geophys.* 45 <https://doi.org/10.1029/2005RG000183>.
- Federal Statistical Office, 2022. Acreage of selected crops in a time comparison [WWW Document]. URL. Federal Statistical Office (accessed 11.17.22). <https://www.destatis.de/EN/Themes/Economic-Sectors-Enterprises/Agriculture-Forestry-Fisheries/Field-Crops-Grassland/Tables/1-acreage-of-selected-crops-in-a-time-comparison.html>.
- Fieuzal, R., Baup, F., Marais-Sicre, C., 2013. Monitoring wheat and rapeseed by using synchronous optical and radar satellite Data—From temporal signatures to crop parameters estimation. *Adv. Remote Sens.* 2, 162–180. <https://doi.org/10.4236/ars.2013.22020>.
- Frantz, D., 2019. FORCE-Landsat + Sentinel-2 analysis ready data and beyond. *Remote Sens.* 11 <https://doi.org/10.3390/rs11091124>.
- Frantz, D., Haß, E., Uhl, A., Stoffels, J., Hill, J., 2018. Improvement of the fmask algorithm for Sentinel-2 images: separating clouds from bright surfaces based on parallax effects. *Remote Sens. Environ.* 215, 471–481. <https://doi.org/10.1016/j.rse.2018.04.046>.
- Gerstmann, H., Doktor, D., Gläßer, C., Möller, M., 2016. PHASE: a geostatistical model for the kriging-based spatial prediction of crop phenology using public phenological and climatological observations. *Comput. Electron. Agric.* 127, 726–738. <https://doi.org/10.1016/j.compag.2016.07.032>.
- Harfenmeister, K., Itzerott, S., Weltzien, C., Spengler, D., 2021. Detecting phenological development of winter wheat and winter barley using time series of Sentinel-1 and Sentinel-2. *Remote Sens.* 13, 5036. <https://doi.org/10.3390/rs13245036>.

- Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. *Neural Comput.* 9, 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>.
- Holtgrave, A.-K., Lobert, F., Erasmi, S., Röder, N., Kleinschmit, B., 2023. Grassland mowing event detection using combined optical, SAR, and weather time series. *Remote Sens. Environ.* 295, 113680 <https://doi.org/10.1016/j.rse.2023.113680>.
- Holtgrave, A.-K., Röder, N., Ackermann, A., Erasmi, S., Kleinschmit, B., 2020. Comparing Sentinel-1 and -2 data and indices for agricultural land use monitoring. *Remote Sens.* 12, 2919. <https://doi.org/10.3390/rs12182919>.
- Huete, A., Didan, K., Miura, T., Rodriguez, E.P., Gao, X., Ferreira, L.G., 2002. Overview of the radiometric and biophysical performance of the MODIS vegetation indices. *Remote Sens. Environ.* 83, 195–213. [https://doi.org/10.1016/S0034-4257\(02\)00096-2](https://doi.org/10.1016/S0034-4257(02)00096-2).
- Jia, M., Tong, L., Zhang, Y., Chen, Y., 2013. Multitemporal radar backscattering measurement of wheat fields using multifrequency (L, S, C, and X) and full-polarization. *Radio Sci.* 48, 471–481. <https://doi.org/10.1002/rds.20048>.
- Jimenez-Perez, G., Alcaine, A., Camara, O., 2019. U-Net Architecture for the Automatic Detection and Delineation of the Electrocardiogram. In: *Computing in Cardiology*. IEEE Computer Society. <https://doi.org/10.23919/CinC49843.2019.9005824>.
- Kaspar, F., Zimmermann, K., Polte-Rudolf, C., 2015. An overview of the phenological observation network and the phenological database of Germany's national meteorological service (Deutscher Wetterdienst). *Adv. Sci. Res.* 11, 93–99. <https://doi.org/10.5194/asr-11-93-2014>.
- Katal, N., Rzanny, M., Mäder, P., Wäldchen, J., 2022. Deep learning in plant phenological research: a systematic literature review. *FrontiersPlant Sci.* 13.
- Kattenborn, T., Leitloff, J., Schiefer, F., Hinz, S., 2021. Review on convolutional neural networks (CNN) in vegetation remote sensing. *ISPRS J. Photogramm. Remote Sens.* 173, 24–49. <https://doi.org/10.1016/j.isprsjprs.2020.12.010>.
- Kattenborn, T., Schiefer, F., Frey, J., Feilhauer, H., Mahecha, M.D., Dormann, C.F., 2022. Spatially autocorrelated training and validation samples inflate performance assessment of convolutional neural networks. *ISPRS Open J. Photogramm. Remote Sens.* 5, 100018 <https://doi.org/10.1016/j.ojphoto.2022.100018>.
- Kavats, O., Khranov, D., Sergieieva, K., Vasylijev, V., 2019. Monitoring harvesting by time series of Sentinel-1 SAR data. *Remote Sens.* 11, 1–16. <https://doi.org/10.3390/rs11212496>.
- Kingma, D.P., Ba, J.L., 2015. Adam: A method for stochastic optimization. *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings* 1–15.
- Kowalski, K., Senf, C., Hostert, P., Pflugmacher, D., 2020. Characterizing spring phenology of temperate broadleaf forests using landsat and Sentinel-2 time series. *Int. J. Appl. Earth Obs. Geoinf.* 92, 102172 <https://doi.org/10.1016/j.jag.2020.102172>.
- Kuhn, M., 2020. *caret: Classification and Regression Training. R package version 6.0-85*.
- Lieth, H. (Ed.), 1974. *Phenology and Seasonality Modeling, Ecological Studies*. Springer-Verlag, New York.
- Lobert, F., 2022. Rcodede: fuctions to use the CODE-DE satellite data repository in R. <https://github.com/felixlobert/rcodede>.
- Lobert, F., Holtgrave, A.-K., Schwieder, M., Pause, M., Vogt, J., Gocht, A., Erasmi, S., 2021. Mowing event detection in permanent grasslands: systematic evaluation of input features from Sentinel-1, Sentinel-2, and landsat 8 time series. *Remote Sens. Environ.* 267, 112751 <https://doi.org/10.1016/j.rse.2021.112751>.
- Löw, J., Ullmann, T., Conrad, C., 2021. The impact of phenological developments on interferometric and polarimetric crop signatures derived from sentinel-1: examples from the DEMMIN study site (Germany). *Remote Sens.* 13, 2951. <https://doi.org/10.3390/rs13152951>.
- Ma, X., Zhu, X., Xie, Q., Jin, J., Zhou, Y., Luo, Y., Liu, Y., Tian, J., Zhao, Y., 2022. Monitoring nature's calendar from space: emerging topics in land surface phenology and associated opportunities for science applications. *Global Change Biology n/a*. <https://doi.org/10.1111/gcb.16436>.
- McNairn, H., Jiao, X., Pacheco, A., Sinha, A., Tan, W., Li, Y., 2018. Estimating canola phenology using synthetic aperture radar. *Remote Sens. Environ.* 219, 196–205. <https://doi.org/10.1016/j.rse.2018.10.012>.
- McNemar, Q., 1947. Note on the sampling error of the difference between correlated proportions or percentages. *Psychometrika* 12, 153–157. <https://doi.org/10.1007/BF02295996>.
- Menzel, A., 2002. Phenology: its importance to the global change community. *Clim. Chang.* 54, 379–385. <https://doi.org/10.1023/A:1016125215496>.
- Menzel, A., Sparks, T.H., Estrella, N., Koch, E., Aasa, A., Ahas, R., Alm-Kübler, K., Bissolli, P., Braslavská, O., Briede, A., Chmielewski, F.M., Crepinsek, Z., Curnel, Y., Dahl, Á., Defila, C., Donnelly, A., Filella, Y., Jatczak, K., Måge, F., Mestre, A., Nordli, Ø., Peñuelas, J., Pirinen, P., Remisová, V., Scheifinger, H., Striz, M., Susnik, A., Van Vliet, A.J.H., Wielgolaski, F.-E., Zach, S., Züst, A.N.A., 2006. European phenological response to climate change matches the warming pattern. *Glob. Chang. Biol.* 12, 1969–1976. <https://doi.org/10.1111/j.1365-2486.2006.01193.x>.
- Mercier, A., Betbeder, J., Baudry, J., Le Roux, V., Spicher, F., Lacoux, J., Roger, D., Hubert-Moy, L., 2020. Evaluation of Sentinel-1 & 2 time series for predicting wheat and rapeseed phenological stages. *ISPRS J. Photogramm. Remote Sens.* 163, 231–256. <https://doi.org/10.1016/j.isprsjprs.2020.03.009>.
- Meroni, M., D'Andrimont, R., Vrieling, A., Fasbender, D., Lemoine, G., Rembold, F., Seguin, L., Verhegghen, A., 2021. Comparing land surface phenology of major European crops as derived from SAR and multispectral data of Sentinel-1 and -2. *Remote Sens. Environ.* 253 <https://doi.org/10.1016/j.rse.2020.112232>.
- Morisette, J.T., Richardson, A.D., Knapp, A.K., Fisher, J.I., Graham, E.A., Abatzoglou, J., Wilson, B.E., Breshears, D.D., Henebry, G.M., Hanes, J.M., Liang, L., 2009. Tracking the rhythm of the seasons in the face of global change: phenological research in the 21st century. *Front. Ecol. Environ.* 7, 253–260. <https://doi.org/10.1890/070217>.
- Nasrallah, A., Baghdadi, N., El Hajj, M., Darwish, T., Belhouchette, H., Faour, G., Darwich, S., Mhawej, M., 2019. Sentinel-1 data for winter wheat phenology monitoring and mapping. *Remote Sens.* 11, 2228. <https://doi.org/10.3390/rs11192228>.
- Pelletier, C., Webb, G.I., Petitjean, F., 2019. Temporal convolutional neural network for the classification of satellite image time series. *Remote Sens.* 11, 1–22. <https://doi.org/10.3390/rs11050523>.
- Perslev, M., Jensen, M.H., Darkner, S., Jennun, P.J., Igel, C., 2019. U-time: a fully convolutional network for time series segmentation applied to sleep staging. *Adv. Neural Inf. Process. Syst.* 32, 1–19.
- Pipia, L., Belda, S., Franch, B., Verrelst, J., 2022. Trends in satellite sensors and image time series processing methods for crop phenology monitoring. In: Bochtis, D.D., Lampridi, M., Petropoulos, G.P., Ampatzidis, Y., Pardalos, P. (Eds.), *Springer Optimization and its Applications*. Springer, Cham, pp. 199–231. https://doi.org/10.1007/978-3-030-84144-7_8.
- R Core Team, 2022. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>.
- Richardson, A.D., Keenan, T.F., Migliavacca, M., Ryu, Y., Sonnentag, O., Toomey, M., 2013. Climate change, phenology, and phenological control of vegetation feedbacks to the climate system. *Agric. For. Meteorol.* 169, 156–173. <https://doi.org/10.1016/j.agrformet.2012.09.012>.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: convolutional networks for biomedical image segmentation. In: *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*. LNCS. Springer, pp. 234–241.
- Scheffler, D., Frantz, D., Segl, K., 2020. Spectral harmonization and red edge prediction of Landsat-8 to Sentinel-2 using land cover optimized multivariate regressors. *Remote Sens. Environ.* 241 <https://doi.org/10.1016/j.rse.2020.111723>.
- Schlund, M., Erasmi, S., 2020. Sentinel-1 time series data for monitoring the phenology of winter wheat. *Remote Sens. Environ.* 246, 111814 <https://doi.org/10.1016/j.rse.2020.111814>.
- Shang, J., Liu, J., Poncos, V., Geng, X., Qian, B., Chen, Q., Dong, T., Macdonald, D., Martin, T., Kovacs, J., Walters, D., 2020. Detection of crop seeding and harvest through analysis of time-series Sentinel-1 interferometric SAR data. *Remote Sens.* 12, 1–18. <https://doi.org/10.3390/rs12101551>.
- Small, D., 2011. Flattening gamma: radiometric terrain correction for SAR imagery. *IEEE Trans. Geosci. Remote Sens.* 49, 3081–3093. <https://doi.org/10.1109/TGRS.2011.2120616>.
- Tetteh, G.O., Gocht, A., Erasmi, S., Schwieder, M., Conrad, C., 2021. Evaluation of Sentinel-1 and Sentinel-2 feature sets for delineating agricultural fields in heterogeneous landscapes. *IEEE Access* 9, 116702–116719. <https://doi.org/10.1109/ACCESS.2021.3105903>.
- Torres, R., Snoeijs, P., Geudtner, D., Bibby, D., Davidson, M., Attema, E., Potin, P., Rommen, B.O., Floury, N., Brown, M., Traver, I.N., Deghaye, P., Duesmann, B., Rosich, B., Miranda, N., Bruno, C., L'Abbate, M., Croci, R., Pietropaolo, A., Huchler, M., Rostan, F., 2012. GMES Sentinel-1 mission. *Remote Sens. Environ.* 120, 9–24. <https://doi.org/10.1016/j.rse.2011.05.028>.
- Veloso, A., Mermoz, S., Bouvet, A., Le Toan, T., Planells, M., Dejoux, J.F., Ceschia, E., 2017. Understanding the temporal behavior of crops using Sentinel-1 and Sentinel-2-like data for agricultural applications. *Remote Sens. Environ.* 199, 415–426. <https://doi.org/10.1016/j.rse.2017.07.015>.
- Vreugdenhil, M., Wagner, W., Bauer-Marschallinger, B., Pfeil, I., Teubner, I., Rüdiger, C., Strauss, P., 2018. Sensitivity of Sentinel-1 backscatter to vegetation dynamics: an austrian case study. *Remote Sens.* 10, 1396. <https://doi.org/10.3390/rs10091396>.
- Wadoux, A.M.J.C., Heuvelink, G.B.M., de Bruin, S., Brus, D.J., 2021. Spatial cross-validation is not the right way to evaluate map accuracy. *Ecol. Model.* 457, 109692 <https://doi.org/10.1016/j.ecolmodel.2021.109692>.
- Ye, Y., Zhang, X., Shen, Y., Wang, J., Crimmins, T., Scheifinger, H., 2022. An optimal method for validating satellite-derived land surface phenology using in-situ observations from national phenology networks. *ISPRS J. Photogramm. Remote Sens.* 194, 74–90. <https://doi.org/10.1016/j.isprsjprs.2022.09.018>.
- Yeasin, M., Haldar, D., Kumar, S., Paul, R.K., Ghosh, S., 2022. Machine learning techniques for phenology assessment of sugarcane using conjunctive SAR and optical data. *Remote Sens.* 14, 3249. <https://doi.org/10.3390/rs14143249>.
- Zeng, L., Wardlow, B.D., Xiang, D., Hu, S., Li, D., 2020. A review of vegetation phenological metrics extraction using time-series, multispectral satellite data. *Remote Sens. Environ.* 237, 111511 <https://doi.org/10.1016/j.rse.2019.111511>.
- Zhang, X., Wang, J., Gao, F., Liu, Y., Schaaf, C., Friedl, M., Yu, Y., Jayavelu, S., Gray, J., Liu, L., Yan, D., Henebry, G.M., 2017. Exploration of scaling effects on coarse resolution land surface phenology. *Remote Sens. Environ.* 190, 318–330. <https://doi.org/10.1016/j.rse.2017.01.001>.
- Zhu, Z., Wang, S., Woodcock, C.E., 2015. Improvement and expansion of the fmask algorithm: cloud, cloud shadow, and snow detection for landsats 4–7, 8, and sentinel 2 images. *Remote Sens. Environ.* 159, 269–277. <https://doi.org/10.1016/j.rse.2014.12.014>.
- Zhu, Z., Woodcock, C.E., 2012. Object-based cloud and cloud shadow detection in landsat imagery. *Remote Sens. Environ.* 118, 83–94. <https://doi.org/10.1016/j.rse.2011.10.028>.