



A Sentinel-2 Machine Learning Dataset for Tree Species Classification in Germany

Maximilian Freudenberg¹, Sebastian Schnell², and Paul Magdon³

¹Chair of Forest Inventory and Remote Sensing & Neural Data Science Group, University of Göttingen, Germany

²Thünen Institute of Forest Ecosystems, Eberswalde, Germany

³Faculty of Resource Management, University of Applied Sciences and Arts (HAWK), Göttingen, Germany

Abstract. We present a machine learning dataset for tree species classification in Sentinel-2 satellite image time series of bottom of atmosphere reflectance. The dataset is based on the German national forest inventory of 2012, as well as analysis ready satellite imagery computed using the FORCE processing pipeline. From the national forest inventory data, we extracted the tree positions, filtered 387 775 trees in the upper canopy layer and automatically extracted the corresponding bottom of atmosphere reflectance time series from Sentinel-2 L2A images. These time series are labeled with the corresponding tree species, which allows pixel-wise classification tasks. Furthermore, we provide auxiliary information such as the approximate tree position, the year of possible disturbance events or the diameter at breast height. Temporally, the dataset spans the years from July 2015 to end of October 2022 with ca. 75.3 million data points for trees of 51 species and species groups, as well as 13.8 million observations for non-tree background. Spatially, it covers entire Germany. The dataset is available under following DOI (Freudenberg et al., 2024): <https://doi.org/10.3220/DATA20240402122351-0>



1 Introduction

In this work, we present a new training dataset for pixel-wise classification of tree species using Sentinel-2 time series of bottom-of-atmosphere (BOA) reflectances across Germany.

Climate change increases the risk of severe weather events such as heavy rainfall or droughts in Central Europe (Toreti et al., 2023). The recent past has seen large scale forest diebacks due to drought, disease or insect manifestations or a combination of these factors (Senf et al., 2020; Senf and Seidl, 2021b). Forest managers face the challenge of adapting their management practices through diversification and other strategies to mitigate these threats. Here, remote sensing will play an increasingly important role as it can support well-informed decisions by providing extensive land cover and forest information at higher temporal frequencies than traditional forest monitoring approaches. In this context, information on tree species is an essential information, key to many forest management decisions.

Tree species classification in satellite imagery is crucial, not only for scientific, but also for practical applications in forestry and nature conservation. This task has been in focus since the early days of space-borne remote sensing with the first Landsat sensors (Walsh, 1980) and it continues today with extensive use of machine-learning methods (Bolyn et al., 2022; Blickeisdörfer et al., 2024).

Sentinel-2 (S2) satellite images are the ideal basis for such analyses, as they are standardized, freely available and collected with high temporal revisit frequency. Machine learning, particularly deep learning, is commonly employed to tackle classification tasks in image data, albeit requiring substantial amounts of training data. In the context of tree species classification, generating training data is demanding and one has to resort to visual interpretation and on-screen labeling of high resolution aerial images, ideally combined with validation in the field – or one has to source labels from forest inventory data.

Ahlswede et al. (2023) have addressed the problem of training data compilation and created a multi-modal training dataset, containing aerial, as well as Sentinel-1 and 2 images of over 50 000 sites in the state of Lower Saxony, Germany. The dataset contains image-wise labels for 20 European tree species, generated from stand level forest inventory data. Utilizing different deep learning models, the authors achieved an F_1 score of 54.6%, using Sentinel-2 data alone. They conclude that “the integration of multi-seasonal data might disentangle further species-related information regarding phenology phases” (Ahlswede et al., 2023, p. 691) – this is what we aim for with the dataset presented here.

Hemmerling et al. (2021) used exactly this kind of multi-seasonal Sentinel-2 data to classify 17 different tree species in the state of Brandenburg, Germany. They applied a random forest classifier to time series of the years 2018 and 2019 and reached F_1 scores between 67% and 99% for the nine most frequent species, thereby demonstrating that at least a subset of species can be separated using S2 time series comparable to the ones provided here. As in the first study, the authors obtained their labels from forest inventories conducted by state authorities.

These two studies are noteworthy exceptions regarding the amount of training data used, because the used datasets were relatively large. Fassnacht et al. (2016) reviewed studies on tree species classification from remotely sensed data and conclude that “investigations focusing on [...] a single often comparably small test site by far dominated the reviewed studies”. This



hinders the generalizability of results and the applicability of generated models to other areas: a dataset covering a large area
45 and long time spans is needed.

To overcome the problem of limited training data we tap the largest dataset of field observations of tree species in Germany:
the national forest inventory (NFI). The German NFI is conducted at full scale every 10 years, with a subsample after 5 years,
and covers more than 25 000 sites, over 60 000 sampling points and more than 500 000 trees across all ownerships and site
50 conditions (Polley et al., 2018). For each tree, several variables such as species, relative position and diameter at breast height
(DBH) are recorded. The resulting dataset is the most comprehensive available for German forests and the derived statistics
provide valuable insights into the forest condition, composition and development on regional and national level. However, the
design of the NFI was not tailored for creating remote sensing reference datasets but to provide an efficient sampling and
plot design for estimating key forest variables. From a remote sensing perspective, one of the major caveats is, that the exact
sampling positions need to be kept confidential, e.g., to prevent biased estimates when management practices are changed in
55 the plot vicinity.

The goal of the work presented is twofold: first, to make satellite data at NFI plot positions available for third parties without
revealing the exact geolocations and second, to analyze the separability and temporal patterns of tree crown reflectances for tree
species in Germany. We link NFI records to BOA reflectance time series from matching Sentinel-2 images, enabling tree species
classification and other applications for a broad range of potential users. Said time series were extracted from analysis ready
60 data generated by the Framework for Operational Radiometric Correction for Environmental monitoring (FORCE) (Frantz,
2019), hosted on the CODE-DE¹ platform. The resulting dataset provides BOA reflectances from July 2015 to October 2022
and in sum contains the time series of 387 775 individual trees and 70,242 non-tree locations. In total there are ca. 75.3
million data points for trees and 13.8 million observations for non-tree background, covering the entirety of Germany and
51 tree species and species groups. The dataset is available online under <https://doi.org/10.3220/DATA20240402122351-0>
65 (Freudenberg et al., 2024) with CC BY 4.0 license.

2 Materials and methods

2.1 Study area and national forest inventory

The dataset covers the entire area of Germany, including islands. More specifically, it contains 24 925 of the 25 382 cluster plots
recorded in the 2012 national forest inventory. Temperate broadleaf and mixed forests prevail in most regions of the country.
70 Coniferous forests, mainly consisting of *Picea abies* (spruce), dominate at higher elevations and forests with *Pinus sylvestris*
(pine) occur on the sandy soil of the north-eastern part of the country. In 2012, about 32% of Germany was covered by forest
(Polley et al., 2018), but due to heavy droughts and following insect infestations in the years 2018–2022 the area of stocked
forest has likely decreased.

¹<https://code-de.org>



The German national forest inventory is conducted on a regular, square sampling grid as shown in Figure 1 with a grid size of $4\text{ km} \times 4\text{ km}$ or less, depending on the federal state. At each grid point there are four inventory plots, aligned in a $150\text{ m} \times 150\text{ m}$ square. The south-western-most inventory plot aligns with the $4\text{ km} \times 4\text{ km}$ grid, as shown in Figure 2.

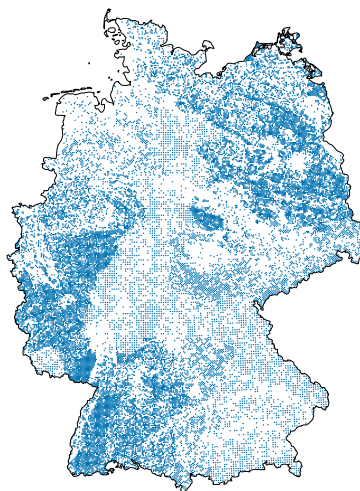


Figure 1. The sampling positions of the German national forest inventory 2012. Borders: © GeoBasis-DE / BKG (2024)

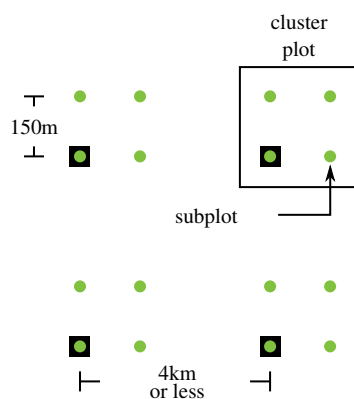


Figure 2. The German national forest inventory sampling grid (black squares) and the inventory points (green). The south-western most inventory point in each cluster plot is aligned with the overarching grid.

The geolocation of each inventory point is measured with a Global Navigation Satellite System (GNSS) device, which may or may not be differentially corrected using correction information from terrestrial reference stations. At this sample point, two angle count samplings are performed (Gregoire and Valentine, 2007), which means that trees whose diameter at breast height (DBH) covers more than a certain solid angle are recorded.



The first angle count sampling includes all trees within a distance from the sample location of 25 times their DBH. For the selected trees, azimuth angle, distance to the plot center, tree species, DBH and other variables are recorded - these measurements form the basis of our labels. A second angle count sampling captures the surrounding forest composition by recording the species of all trees within a radius of 50 times their DBH - it samples trees up to larger distances compared to the first
85 sampling. The second angle count sampling allows to tell, which sub-plots are pure stands, i.e. have only one tree species in them. This in turn allows to mark a subset of tree species labels with high confidence because they grow in stands that are most likely composed of only one species; the information whether a stand is pure is included in the dataset.

2.2 NFI reference data selection

To compile the provided training dataset we used the NFI data in the following way: First, we removed all trees that grow in the
90 understory; this information is recorded during the inventory. For the remaining trees we modeled a circular stand area using species specific parameters as provided in (Riedel et al., 2017, pp. 39, 40). As we know the position of each tree, as well as its estimated stand area, we can remove trees that are probably not visible from above by a heuristic.

We count trees as visible when they are either the biggest (area-wise) within a radius of 3 m or there are no other trees within that radius. Furthermore, we count them as visible, if their stand area overlaps the union of all other stand areas by not more
95 than 50%, as depicted in Figure 3. Trees labeled as visible by this heuristic form the basis for the training dataset.

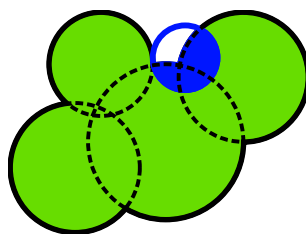


Figure 3. Sketch of a tree group: Green trees are assumed to be visible. The blue tree overlaps with more than 50% of its area with other trees and is therefore discarded.

To allow training classification methods for the discrimination between tree and non-tree pixels we added non-forest observations to the dataset. For this, we sampled the tree cover density layer provided by the Copernicus Land Monitoring Service for the year 2018 in the vicinity of the NFI plots². The tree cover density layer is sampled at locations that are at least 20 meters away from the next pixel with tree density greater than 10%.

100 2.3 Satellite data selection

We used images from the Sentinel-2 satellites, pre-processed to analysis-ready level by the FORCE processing pipeline (Frantz, 2019). FORCE provides a way to compute harmonized time series that are spatially and spectrally well aligned, which is discussed in more detail later. The resulting data comprises all S2 bands with 10 or 20 m resolution. Additionally, FORCE

²<https://land.copernicus.eu/en/products/high-resolution-layer-tree-cover-density>



provides quality assurance information (QAI) that aids in filtering out undesirable image conditions such as clouds, snow, or high water vapor content. The data is hosted on the CODE-DE³ and EO-Lab⁴ platforms. End users have the option to either download the pre-processed data or can re-process it using the same settings utilized in generating the FORCE data cube on CODE-DE. The necessary parameter files are provided alongside the dataset.

2.4 Time series extraction and data processing

From the FORCE data cube we clipped $300\text{ m} \times 300\text{ m}$ image patches containing the 24 925 filtered NFI cluster plots and their surroundings, as depicted in Figure 4. We extracted the bottom of atmosphere reflection (BOA) as well as the quality assurance information (QAI). Before extraction, we filtered the plots to ensure they contained at least one pixel with data, not affected by clouds or cloud shadows.

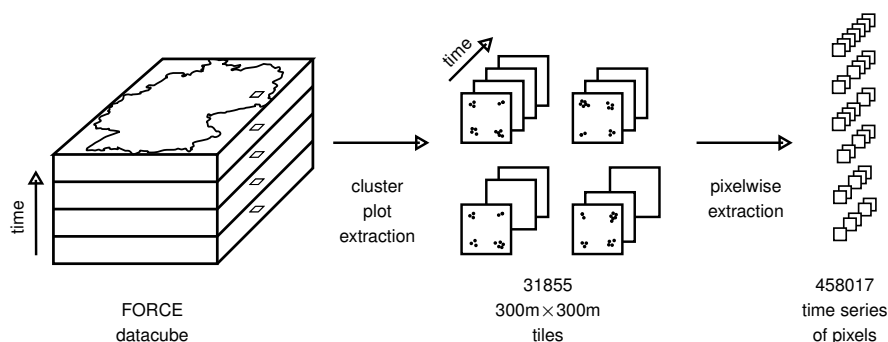


Figure 4. The time series extraction workflow: First, $300\text{ m} \times 300\text{ m}$ tiles are clipped from the FORCE datacube for Germany for all records between 2012 and 2022. Second, the pixel-wise time series are extracted from the tile time series.

In a last step, we extracted the BOA and QAI pixel time series from the extracted patches at the respective reference data positions. In cases where a single tree covered more than one $10\text{ m} \times 10\text{ m}$ Sentinel pixel, we calculated the area-weighted average of all pixels intersected by the tree's crown projection area, as depicted in Figure 5. Each extracted satellite observation was then linked to its acquisition date, the corresponding NFI data and more information. Senf and Seidl (2021a) provide a Landsat-based map of forest disturbances for Germany between 1986 and 2020 at a resolution of 30 m . To be able to identify possible disturbance events, we included the disturbance year from this map in the dataset. However, this still leaves a gap between 2020 and 2022, for which no disturbance information is available. We bridged this gap by attaching the information whether the trees were still present during the 2022 NFI. To enable approximate spatial analyses, we furthermore included the center coordinate of the 1 km Inspire-grid tile the cluster plots are located in.

The final dataset comprises the following data and an excerpt is given in Table A1 in the appendix:

³<https://code-de.org>

⁴<https://eo-lab.org>

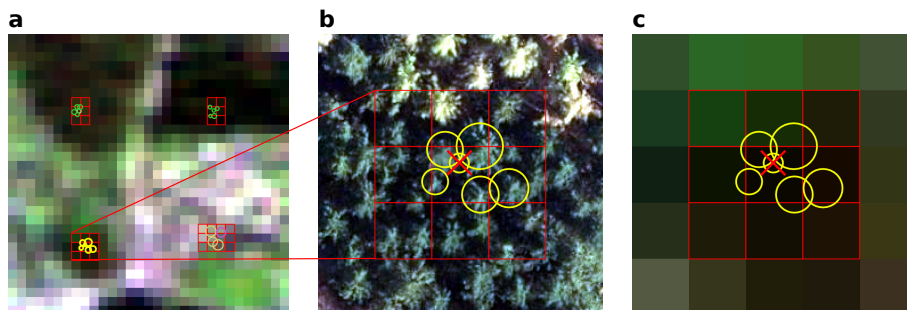


Figure 5. (a) The whole cluster plot cutout of 300 m×300 m. S2-Image: European Space Agency (2021) (b) The lower left subplot with the corresponding orthophoto for reference. Douglas firs in the lower part, spruce in the upper part of the image. Image: © BKG (2021) (c) The S2 pixels corresponding to the subplot with circles depicting the modeled tree crown areas. The crossed out tree is omitted because it overlaps too much with surrounding trees.

- Identifiers for individual trees: a global ID, the NFI cluster plot ID (tnr) and the corner ID (enr) within the plot. Non-forest records have negative IDs.
- 125 – The tree species encoded according to the official NFI schema, provided within the dataset in a separate table "x_ba".
- The acquisition date, encoded as Unix time, representing the number of seconds elapsed since January 1, 1970, 00:00 UTC. Every date was randomly shifted by up to three days.
- The BOA reflectance values: 10 signed 16-bit integers, one for each band, encoded as 20 byte blob data. To hamper the identification of exact plot positions, each value was multiplied with a uniform random number between 0.95 and 1.05.
- 130 – Quality assurance information bit-flags, encoded as 16-bit integers, allowing for filtering based on image quality. The FORCE documentation provides details on the meaning of each bit⁵.
- Diameter at breast height in millimeters, tree height in decimeters, modeled according to Riedel et al. (2017) and crown area in m² allowing for further tree filtering or analyses by diameter class.
- The WGS84 center coordinate of the 1 km Inspire grid tile the tract can be found in.
- 135 – The disturbance year according to the map provided by Senf and Seidl (2021a).
- Whether the NFI position measurement was differentially corrected.
- Whether the record belongs to the train or validation set (see below).
- Whether the record comes from a pure stand according to the class definitions of the NFI.

⁵<https://force-eo.readthedocs.io/en/latest/howto/qai.html#quality-bits-in-force>



- Whether the tree was observed again in the 2022 forest inventory.
- 140 – The day of year of the acquisition, corresponding to the shifted date.

All samples were randomly split into training and validation sets based on their cluster plot IDs with a ratio of 70% - 30%. This rules out any spatial overlap between the training and test sets and reduces correlations between the two. For benchmark studies, we recommend using this split to ensure comparability across publications.

2.5 Assessment of the geolocation accuracy of the NFI plots

- 145 The tree positions in the NFI are measured in polar coordinates relative to the plot center, using a compass for the angle and an ultrasonic device for the distance measurement. We assume that the errors for angle and distance are small compared to the GNSS error of the plot center position measurement. GNSS measurements can be differentially corrected by using ground-based reference stations to increase positional accuracy. Depending on the federal state and field team, coordinates of the plot centers are measured with corrected GNSS devices or not. Of the sub-plots with trees in the dataset, 76.5% were corrected,
- 150 22.5% were not, and the remainder has unknown status.

To estimate the accuracy of the plot center coordinates, we compared the field-measured tree positions with tree positions derived from true-ortho aerial images, obtained from the Federal Agency for Cartography and Geodesy. These images are ortho-rectified using a surface model and aligned with high accuracy to ground control points. The ATKIS orthophoto standard guarantees a geolocation error with standard deviation of 0.4 m or less (Arbeitsgemeinschaft der Vermessungsverwaltungen der Länder der Bundesrepublik Deutschland (AdV), 2020). Two expert image interpreters then manually shifted a sample of

155 200 NFI plot positions, and thereby the trees, to match the true tree positions by comparing local tree patterns as depicted in Figure 6. This allows to quantitatively evaluate the deviation of measured from true positions and to compare the accuracy of corrected and uncorrected measurements.

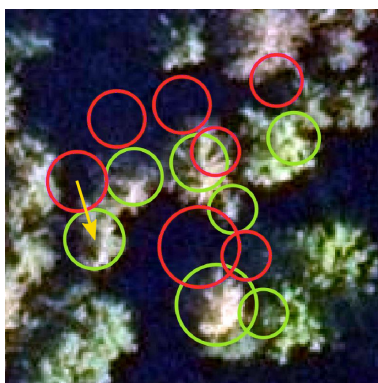


Figure 6. Original, measured GNSS coordinates (red) were shifted by 4.8 m to the visually best matching position (green) in aerial orthophotos to quantify GNSS errors. Circles depict modeled stand areas.



3 Dataset description and statistics

160 3.1 Numerical species distribution

Due to the highly varying dominance of tree species in Germany, the numerical distribution of the different species (Figure 7) is heavily imbalanced. The most abundant species is *Pinus sylvestris* (Pine), followed by *Picea abies* (Spruce), *Fagus sylvatica* (Beech) and the different *Quercus* (Oak) species. Note that all statistics only represent the dataset used here and not the NFI itself, albeit both are closely related. For a list of all included tree species and their counts we refer to appendix Table A3.

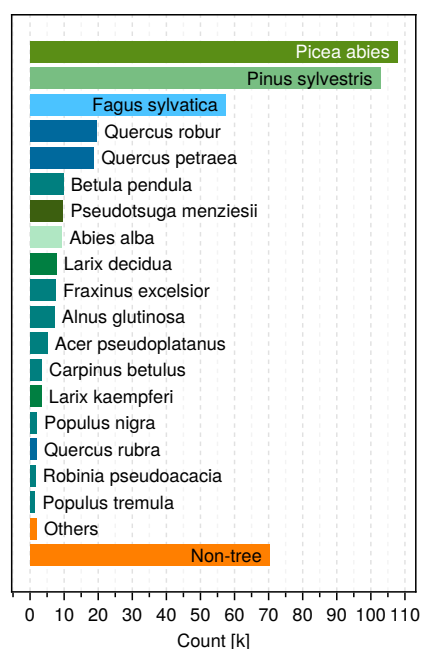


Figure 7. The numerical species distribution in the training dataset for all trees extracted from the NFI.

165 3.2 Temporal signatures of selected species

Coniferous and deciduous trees can be clearly separated visually by inspecting the time series of their infrared (IR) reflectance, as depicted in Figure 8. In the presented time series, the observations for a given species and point in time have been averaged across all undisturbed individuals in pure stands. Whether a stand is pure or not was determined using the second angle count sampling of the NFI. Obviously, deciduous trees exhibit a much stronger seasonal pattern than coniferous trees. However, this separation is less evident in the green band a) due to its higher susceptibility to atmospheric effects and b) due to its lower absolute reflectance, which deteriorates the signal to noise ratio. While the temporal infrared profiles of *Fagus sylvatica* and *Quercus robur* are generally distinguishable across most years, there are instances where differentiation becomes challenging (e.g. 2016 and 2020). *Quercus robur* tends to have a slightly lower IR reflectance on average, particularly in summer. *Picea*



175 abies and *Pinus sylvestris* also differ only slightly in the infrared, with *Picea abies* having lower average values on trend. Overall, differentiating species by their temporal profiles alone seems challenging without considering their spectrum at the same time.

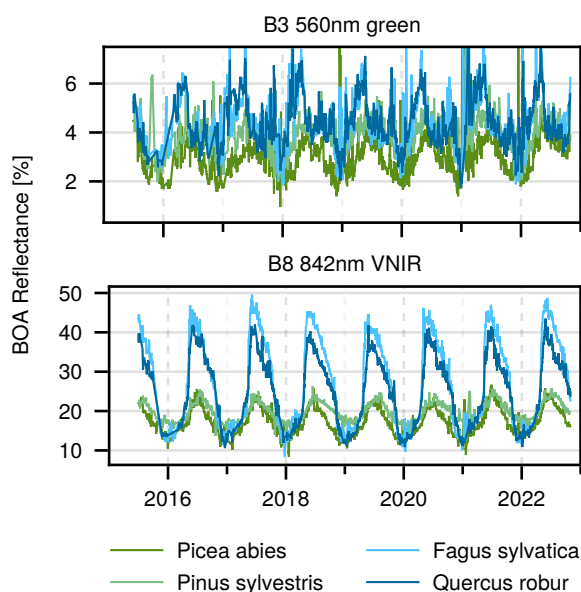


Figure 8. Time series of BOA reflectance for indicated species, averaged over all undisturbed individual trees in pure stands at a given time. The data has been filtered to exclude all types of cloud cover and their shadows, snow, and pixels with high aerosol optical depth.

Looking at a random selection of four individual trees' time series, depicted in Figure 9, it becomes clear that at a single tree level the differences between species seem to be still present, but with high variance from year to year.

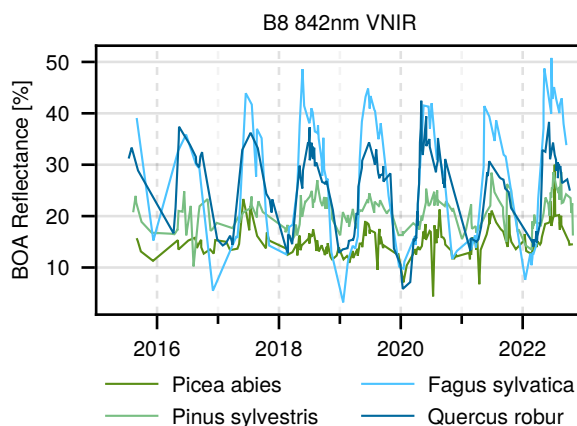


Figure 9. Time series of random single trees of different species.



180 Figure 10 shows the total observation count over time. After the commissioning of Sentinel-2B in June 2017 the number of observations increases. As one would expect, there are more observations in the summer months when clouds are less likely and especially from 2018 onward the counts regularly reach over 1 million.

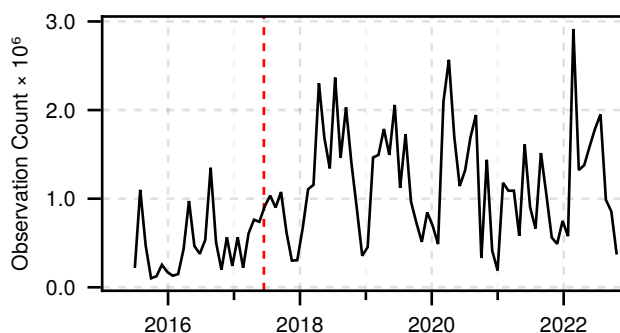


Figure 10. Total monthly observations of all selected pixels in the training dataset. The vertical red line corresponds to the Sentinel-2B commissioning date.

3.3 Spectral signatures

Besides the temporal variation of the reflectance, the spectral variation is an important feature for the tree species classification – however, the species are not necessarily separable by their spectrum alone, as can be seen in Figure 11. It depicts the Sentinel-2 spectra of the five most frequent species, as well as the background class. *Fagus sylvatica* and *Quercus petraea* for example have almost matching spectra, especially in the shorter wavelengths. The resulting spectra match the ones presented in Immitzer et al. (2016).

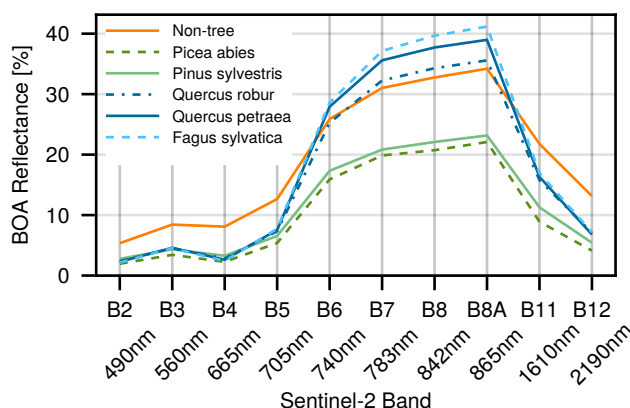


Figure 11. Average spectrum of the five most frequent species in the dataset, plus the background class. Records of pure stands have been averaged between May and August (inclusive) of the years 2017–2022.



3.4 Spatial distribution

It can be expected that the temporal signatures vary with local conditions, e.g. along an latitudinal or elevation gradient. Therefore, it is important to analyze the spatial coverage of the training data. Figure 12 shows that *Picea abies* is mainly present in the south-west of Germany and in the lower mountain ranges. *Pinus sylvestris* on the other hand, is predominant on the sandy soils of the north-eastern part of the country. The different *Quercus* species occur mostly in the west of Germany, but are also present throughout the rest of the country. *Fagus sylvatica*, lastly, co-occurs with *Quercus* sp., but in contrast to them, manages to settle in the higher and therefore colder hillsides of the central parts of Germany. Note however, that these spatial distributions are derived from the dataset, which does not mirror the NFI one to one due to filtering and the availability of satellite images.

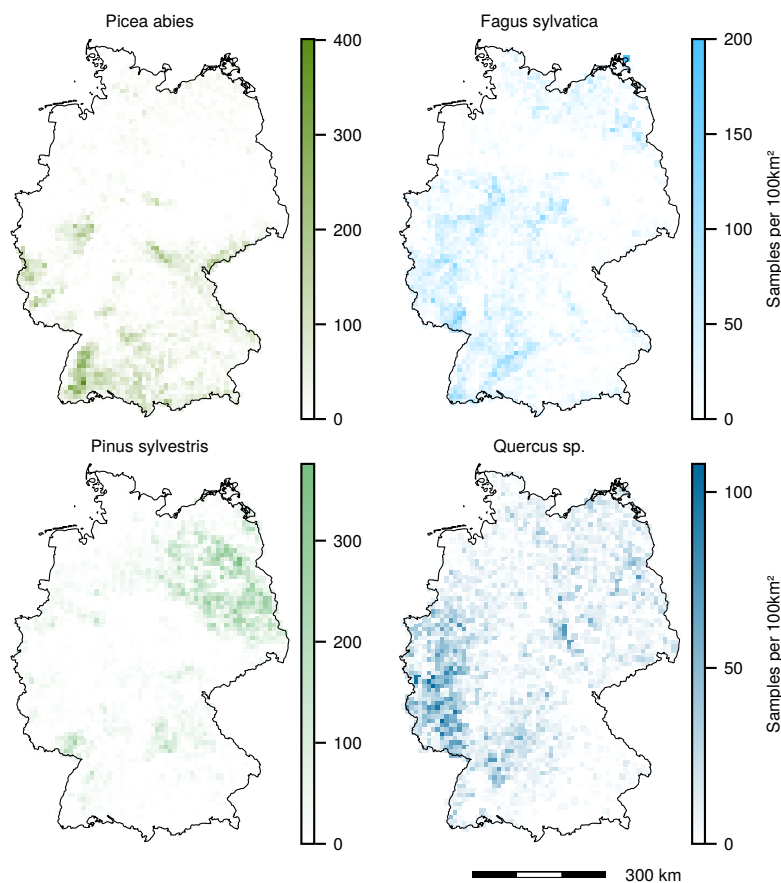


Figure 12. Spatial tree distribution for different species. Note the different scales. Borders: © GeoBasis-DE / BKG (2024)



3.5 NFI geolocation accuracy estimation

The analysis of the spatial accuracy of the NFI plot coordinate GNSS measurements reveals that ninety-five percent of the measured deviations of corrected GNSS positions were smaller than 11.2 m, and 81% were smaller than 5 m; Figure 13 depicts the corresponding histogram along with the empirical cumulative density function. Against expectations, the comparison of corrected and uncorrected GNSS measurements shows no significant difference.

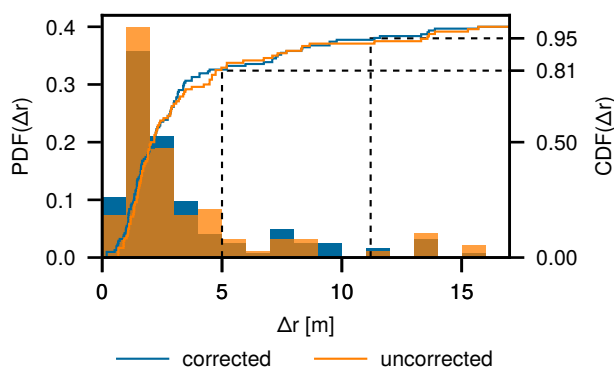


Figure 13. Histogram of distances by which plot locations were shifted from the original GNSS positions. Differentially corrected measurements are depicted in blue.

4 Discussion

4.1 Geolocation accuracy

Sentinel-2: To obtain the presented dataset, we linked spatial information from two different data sources: georeferenced satellite images and on-ground GNSS measurements. A misalignment of these sources might lead to labeling errors in the dataset. FORCE co-registers all Sentinel-2 images with averaged Landsat time series. The Landsat images are in turn co-registered with the Sentinel-2 global reference image which results in a geometric accuracy of 10.2 m at the 90% confidence level for Landsat 8 (Haque et al., 2022). Consequently, this is the best estimate for the spatial accuracy of the used S2 images. The reason for this cyclic co-registration of Sentinel to Landsat to Sentinel is, that so far only the S2 level 1 archive has been processed to a common standard⁶. The level 2 data, which compensates atmospheric effects and is needed for coherent time series, is not yet available at a standardized processing baseline in any public archive.

NFI geolocation accuracy: The comparison of corrected and uncorrected GNSS measurements showed no significant difference in spatial accuracy, at least not the way we measured it. As differential correction unquestionably increases the GNSS accuracy, we suppose that increasing the count of sampled plots as well as the number of image interpreters would change our result. Either way, as the satellite image resolution is 10 m and 81% of the GNSS measurements had an error of less than 5 m,

⁶<https://sentinels.copernicus.eu/web/sentinel/technical-guides/sentinel-2-msi/copernicus-sentinel-2-collection-1-availability-status>



we are confident that the GNSS positions can be combined with the satellite data. It will be interesting to analyze the accuracy of trained classifiers as a function of correction status.

4.2 Adverse imaging conditions

220 During the extraction process, we filtered out most pixels with cloud cover or cloud shadows. FORCE employs the FMASK algorithm (Zhu and Woodcock, 2012) for cloud detection, which has an accuracy of 84% for cloud / clear detection and 72% detection accuracy for cloud shadows (Aybar et al., 2022). Consequently, falsely labeled image regions lead to commission or omission errors in the final dataset, i.e. usable pixels might have been removed by being labeled as cloudy or cloud pixels could be in the dataset. However, there are other imaging conditions that might affect the quality of a pixel like high aerosol content, snow or poor illumination conditions. FORCE encodes this information in the quality assurance information and end
225 users can use this to further narrow the dataset down to only the highest quality pixels.

4.3 Taxonomic identification

The field teams of the NFI data are trained and undergo testing before being allowed to take samples. However, it cannot be ruled out that under adverse conditions certain species are confused. We cannot quantify this error, but assume that the vast majority of tree species identifications are correct, in particular for the common species.

230 4.4 Mixed pixels

At present, we cannot quantify the effect of pixels that contain different tree species on our dataset, as it is in most cases impossible to derive the species shares of a pixel based on the NFI data. The NFI does not fully sample a given plot, so in most cases, labels are only available for parts of a given pixel. Another source for mixed pixels are the 20 m resolution bands of Sentinel-2 that are resampled to 10 m by FORCE, thereby distributing identical information across several pixels.

235 4.5 Species separability analysis

To detect inconsistencies within the dataset, we computed the infrared reflectance histograms of five species for mixed and pure stands. If the histogram shows artifacts like double peaks or differs strongly between pure and mixed stands, this could hint to deficiencies in the respective part of the dataset. Figure 14 shows the histograms of S2 band B8 (842nm) averaged over all records in June 2021 for species whose occurrence is correlated – *Betula pendula* often grows along with *Pinus sylvestris* and *Fagus sylvatica* and *Quercus* spp. often appear together. June 2021 has been chosen because both Sentinel satellites were
240 operational and, unlike the preceding years and 2022, 2021 was not particularly dry.

The reflectance distributions for *Pinus* and *Betula* clearly differ between mixed and pure stands. In mixed stands the distributions are relatively wide and overlap, whereas there are separable peaks for pure stands (albeit some overlap remains) and the distance between maxima is larger. We interpret this as a hint that the dataset contains false labels due to insufficient

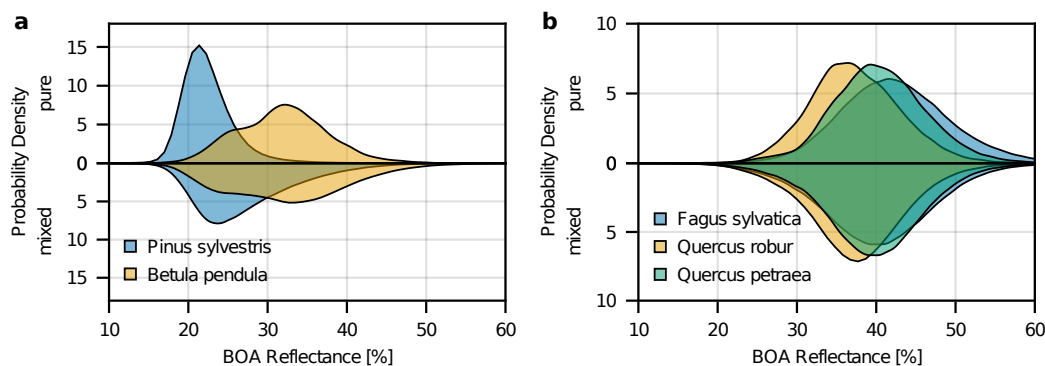


Figure 14. Histogram of near infrared (842 nm) BOA reflectances, averaged over all trees in June 2021, for (a) *Pinus sylvestris* and *Betula pendula* and (b) *Fagus sylvatica*, *Quercus robur* and *Quercus petraea*. The upper parts depict pure stands and the lower parts mixed stands.

245 spatial accuracy or that the extracted pixel values come from mixed pixels containing other species or even different land cover classes.

Comparing *Fagus sylvatica* to two *Quercus* species, one can see that the distributions overlap much more, as all three species are deciduous. In mixed stands there is hardly any observable difference between the distributions. For pure stands the distributions still overlap significantly, but the distance between peaks is slightly larger than for mixed stands. However, the
250 overlap of these distributions does not necessarily indicate labeling errors; it could also be that these are the naturally occurring values. This highlights the necessity of factors beyond spectral data, e.g. temporal profiles as shown in Figure 9, for species classification.

5 Conclusion and outlook

In this work we presented the so far most comprehensive dataset of annotated Sentinel-2 time series data for tree species
255 detection in Germany. With over 380 thousand trees of 48 species observed for over seven years, this dataset can significantly advance research into automatic tree species classification for Germany, and central Europe. At the same time the described approach can serve as a pilot study for making national forest inventory data from other countries accessible for the remote sensing community e.g. for training machine learning models without releasing the exact geolocations publicly. Lessons learned from its application can be used to enhance future inventories and datasets. For example, it could show that for underrepresented
260 species more labels are required, which in turn could be sampled in targeted inventories.

As discussed in the previous section, the dataset still has several shortcomings that could be improved. To achieve better agreement between labels and images, the spatial accuracy of the data sources has to be increased. To do so, we suggest that in future all NFI position measurements are taken using differential GNSS devices, although we saw no significant differences in accuracy. Furthermore, we expect that aligning the Sentinel-2 images directly with the S2 global reference image instead of



265 averaged Landsat time series would improve their spatial accuracy and make it easier to derive interpretable error metrics. We
consider releasing an updated dataset version as soon as Sentinel-2 L2A collection one is fully accessible.

The main focus of further efforts will be to increase the number of labels for weakly represented classes, e.g. by utilizing
automatically classified high resolution orthophotos as reference. First attempts to automatically identify underrepresented tree
species in standard RGBI aerial images with 20 cm spatial resolution have failed, so the presented dataset is still limited regard-
270 ing less abundant species. Another option to increase the overall amount of data would be to incorporate forest inventory data
at the stand level from e.g. state forest enterprises, however, this data often only provides estimates of tree species proportions
within management units, but no geolocation.

We hope that this dataset fosters the research into time series based classification of tree species and believe it offers many
possibilities for analyses that go beyond the ones presented here. Using classification methods in general, one could investigate
275 which spectral bands and which points in time are crucial for precise species classification. As the dataset not only contains
the time series of individual trees' BOA reflectances, but also their approximate location, spatio-temporal patterns in tree
phenology could be assessed on individual species level. For example, the onset of leaf emergence could be analyzed first
in the dataset alone, and later by using species maps generated by a derived classification method. Lastly, the dataset could
be used to correlate reflectances and approximate health conditions with meteorological events like droughts on a per-species
280 level. This would open up further research into climate-change resistant species and enables the identification of endangered
forest stands. In the future we plan to release updated versions of the dataset, particularly after the final publication of the 2022
NFI.

6 Data availability

All data is available online under <https://doi.org/10.3220/DATA20240402122351-0> (Freudenberg et al., 2024) with CC BY 4.0
285 license.

Author contributions. MF: coding, assembling the dataset, main work on manuscript, SS: data provision, proof reading, advice in research
questions, PM: advice in research questions, manuscript development, proof reading

Competing interests. The authors declare to have no conflicts of interest.

Acknowledgements. The authors thank the Thünen Institute of Forest Ecosystems for providing the national forest inventory data. MF thanks
290 Alexander Ecker for ongoing financial support and reviews, as well as Christoph Kleinn for proof-reading.



Financial support. The *Klimba* project and this work were funded by the Federal Ministry for Digital and Transport under grant number 50EW2012A/B.

Appendix A: Database excerpt and species counts

Table A1. Database excerpt. The bottom of atmosphere (BOA) reflectance is encoded as 10 signed 16 bit integers, the quality assurance information (QAI) is a single 16 bit integer. DOY abbreviates day of year.

cluster ID (tnr)	corner ID (enr)	tree ID	species	unix time	BOA	QAI	train	pure	
455	1	69831	211	1440374400	10 16-bit integers	8192	1	0	...
455	1	69831	211	1448064000	10 16-bit integers	10256	1	0	...
455	1	69831	211	1455494400	10 16-bit integers	10240	1	0	...
455	1	69831	211	1460592000	10 16-bit integers	8192	1	0	...
455	1	69831	211	1463961600	10 16-bit integers	8192	1	0	...
455	1	69831	211	1467072000	10 16-bit integers	8192	1	0	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

	DBH [mm]	height [dm]	area [m ²]	X	Y	corrected	disturbance year	present 2022	DOY
...	231	243	20.4	9.80714	47.64294	1	0	1	236
...	231	243	20.4	9.80714	47.64294	1	0	1	325
...	231	243	20.4	9.80714	47.64294	1	0	1	46
...	231	243	20.4	9.80714	47.64294	1	0	1	105
...	231	243	20.4	9.80714	47.64294	1	0	1	144
...	231	243	20.4	9.80714	47.64294	1	0	1	180
...	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮



Table A2. List of all included tree species with counts (part 1).

species code	species	common name	count
-1	-	other land cover	70242
10	Picea abies	Norway spruce	107798
12	Picea sitchensis	sitka spruce	937
19	Picea spec.	other spruces	232
20	Pinus sylvestris	Scots pine	102730
21	Pinus mugo	mountain pine	88
22	Pinus nigra	European black pine	606
24	Pinus cembra	Swiss pine	3
25	Pinus strobus	eastern white pine	431
29	Pinus spec.	other pines	65
30	Abies alba	silver fir	9375
33	Abies grandis	grand fir	384
39	Abies spec.	other firs	291
40	Pseudotsuga menziesii	Douglas fir	9598
50	Larix decidua	European larch	7674
51	Larix kaempferi	Japanese larch (+hybrids)	3308
90		other coniferous trees	139
94	Taxus baccata	European yew	11
100	Fagus sylvatica	beechn	57341
110	Quercus robur	English oak	19617
111	Quercus petraea	sessile oak	18697
112	Quercus rubra	Northern red oak	1861
120	Fraxinus excelsior	common ash	7469
130	Carpinus betulus	hornbeam	3411
140	Acer pseudoplatanus	sycamore maple	5042
141	Acer platanoides	Norway maple	598
142	Acer campestre	field maple	387
150	Tilia spec.	linden tree (indigenous species)	1294
160	Robinia pseudoacacia	black locust	1553
170	Ulmus spec.	elm, native species	406
181	Castanea sativa	chestnut	416
190		misc. deciduous trees with long life expectancy	246
191	Sorbus domestica	service tree	2
193	Sorbus aria	common whitebeam	51
200	Betula pendula	silver birch	9729
201	Betula pubescens	moor birch	858
211	Alnus glutinosa	black alder	7098
212	Alnus incana	grey alder	460
220	Populus tremula	common aspen	1402
221	Populus nigra	European black poplar (+ hybrids)	1945



Table A3. List of all included tree species with counts (part 2).

species code	species	common name	count
222	Populus x canescens	grey poplar (+hybrids)	196
223	Populus alba	silver poplar	109
224	Populus trichocarpa x maximoviczii	balsam poplar	636
230	Sorbus aucuparia	European rowan	270
240	Salix spec.	willow	1203
250	Prunus padus	bird cherry	77
251	Prunus avium	wild cherry	1357
252	Prunus serotina	black cherry	132
290		misc. deciduous trees with short life expectancy	92
292	Malus sylvestris	European crab apple	37
293	Pyrus communis	European wild pear	42
295	Sorbus torminalis	wild service tree	71

References

- 295 Ahlswede, S., Schulz, C., Gava, C., Helber, P., Bischke, B., Förster, M., Arias, F., Hees, J., Demir, B., and Kleinschmit, B.: *TreeSatAI Benchmark Archive*: a multi-sensor, multi-label dataset for tree species classification in remote sensing, *Earth System Science Data*, 15, 681–695, <https://doi.org/10.5194/essd-15-681-2023>, publisher: Copernicus GmbH, 2023.
- Arbeitsgemeinschaft der Vermessungsverwaltungen der Länder der Bundesrepublik Deutschland (AdV): Produkt- und Qualitätsstandard für Digitale Orthophotos, Tech. rep., 2020.
- 300 Aybar, C., Ysuhuaylas, L., Loja, J., Gonzales, K., Herrera, F., Bautista, L., Yali, R., Flores, A., Diaz, L., Cuenca, N., Espinoza, W., Prudencio, F., Llactayo, V., Montero, D., Sudmanns, M., Tiede, D., Mateo-García, G., and Gómez-Chova, L.: CloudSEN12, a global dataset for semantic understanding of cloud and cloud shadow in Sentinel-2, *Scientific Data*, 9, 782, <https://doi.org/10.1038/s41597-022-01878-2>, 2022.
- Blickensdörfer, L., Oehmichen, K., Pflugmacher, D., Kleinschmit, B., and Hostert, P.: National tree species mapping using Sentinel-1/2 time series and German National Forest Inventory data, *Remote Sensing of Environment*, 304, 114 069, 2024.
- 305 Bolyn, C., Lejeune, P., Michez, A., and Latte, N.: Mapping tree species proportions from satellite imagery using spectral–spatial deep learning, *Remote Sensing of Environment*, 280, 113 205, <https://doi.org/https://doi.org/10.1016/j.rse.2022.113205>, 2022.
- European Space Agency: Copernicus Sentinel-2 (processed by ESA), 2021, MSI Level-1C TOA Reflectance Product. Collection 1, Available online: https://doi.org/10.5270/S2_-742ikth, 2021.
- 310 Fassnacht, F. E., Latifi, H., Stereńczak, K., Modzelewska, A., Lefsky, M., Waser, L. T., Straub, C., and Ghosh, A.: Review of studies on tree species classification from remotely sensed data, *Remote Sensing of Environment*, 186, 64–87, <https://doi.org/10.1016/j.rse.2016.08.013>, 2016.
- Frantz, D.: FORCE—Landsat + Sentinel-2 Analysis Ready Data and Beyond, *Remote Sensing*, 11, 1124, <https://doi.org/10.3390/rs11091124>, 2019.
- 315 Freudenberg, M., Schnell, S., and Magdon, P.: Sentinel-2 machine learning dataset for tree species classification in Germany, <https://doi.org/https://doi.org/10.3220/DATA20240402122351-0>, 2024.
- Gregoire, T. G. and Valentine, H. T.: Sampling strategies for natural resources and the environment, CRC Press, 2007.



- Haque, M. O., Rengarajan, R., Lubke, M., Hasan, M. N., Shrestha, A., Tuli, F. T. Z., Shaw, J. L., Denevan, A., Franks, S., Micijevic, E., Choate, M. J., Anderson, C., Thome, K., Kaita, E., Barsi, J., Levy, R., and Miller, J.: ECCOE Landsat Quarterly Calibration and Validation Report—Quarter 3, 2022, <https://pubs.usgs.gov/of/2023/1013/ofr20231013.pdf>, 2022.
- 320 Hemmerling, J., Pflugmacher, D., and Hostert, P.: Mapping temperate forest tree species using dense Sentinel-2 time series, *Remote Sensing of Environment*, 267, 112 743, <https://doi.org/10.1016/j.rse.2021.112743>, 2021.
- Immitzer, M., Vuolo, F., and Atzberger, C.: First Experience with Sentinel-2 Data for Crop and Tree Species Classifications in Central Europe, *Remote Sensing*, 8, 166, <https://doi.org/10.3390/rs8030166>, number: 3 Publisher: Multidisciplinary Digital Publishing Institute, 325 2016.
- Polley, H., Hennig, P., Kroither, F., Marks, A., Riedel, T., Schmidt, U., Schwitzgebel, F., and Stauber, T.: *Der Wald in Deutschland*, Bundesministerium für Ernährung und Landwirtschaft, Wilhelmstraße 54, 10117 Berlin, 3rd, corrected edn., www.bmel.de/publikationen, 2018.
- Riedel, T., Hennig, P., Kroither, F., Polley, H., Schmitz, F., and Schitzgebel, F.: *Die dritte Bundeswaldinventur: BWI 2012; Inventur- und Auswertungsmethoden*, publisher: TI: Johann Heinrich von Thünen-Institut, 2017.
- 330 Senf, C. and Seidl, R.: Mapping the forest disturbance regimes of Europe, *Nature Sustainability*, 4, 63–70, <https://doi.org/10.1038/s41893-020-00609-y>, number: 1 Publisher: Nature Publishing Group, 2021a.
- Senf, C. and Seidl, R.: Persistent impacts of the 2018 drought on forest disturbance regimes in Europe, *Biogeosciences*, 18, 5223–5230, <https://doi.org/10.5194/bg-18-5223-2021>, 2021b.
- 335 Senf, C., Buras, A., Zang, C. S., Rammig, A., and Seidl, R.: Excess forest mortality is consistently linked to drought across Europe, *Nature Communications*, 11, 6200, <https://doi.org/10.1038/s41467-020-19924-1>, number: 1 Publisher: Nature Publishing Group, 2020.
- Toreti, A., Bavera, D., Acosta Navarro, J., Arias Muñoz, C., Barbosa P., De Jager, A., Di Ciollo, C., Fioravanti, G., Hrast Essenfelder, A., Maetens, W., Magni, D., Masante, D., Mazzeschi, M., McCormick, N., and Salamon, P.: *Drought in Europe: August 2023 : GDO analytical report*, Publications Office of the European Union, Luxembourg, ISBN 978-92-68-07670-5, oCLC: 1404455500, 2023.
- 340 Walsh, S. J.: Coniferous tree species mapping using LANDSAT data, *Remote Sensing of Environment*, 9, 11–26, [https://doi.org/https://doi.org/10.1016/0034-4257\(80\)90044-9](https://doi.org/https://doi.org/10.1016/0034-4257(80)90044-9), 1980.
- Zhu, Z. and Woodcock, C. E.: Object-based cloud and cloud shadow detection in Landsat imagery, *Remote Sensing of Environment*, 118, 83–94, <https://doi.org/10.1016/j.rse.2011.10.028>, 2012.