

Clark University

## Clark Digital Commons

---

Geography

Faculty Works by Department and/or School

---

2020

### Accounting for training data error in machine learning applied to earth observations

Arthur Elmes  
*Clark University*

Hamed Alemohammad  
*Radiant Earth Foundation*

Ryan Avery  
*University of California, Santa Barbara*

Kelly Caylor  
*University of California, Santa Barbara*

J. Ronald Eastman  
*Clark University*

*See next page for additional authors*

Follow this and additional works at: [https://commons.clarku.edu/faculty\\_geography](https://commons.clarku.edu/faculty_geography)



Part of the [Remote Sensing Commons](#)

---

#### Repository Citation

Elmes, Arthur; Alemohammad, Hamed; Avery, Ryan; Caylor, Kelly; Eastman, J. Ronald; Fishgold, Lewis; Friedl, Mark A.; Jain, Meha; Kohli, Divyani; Bayas, Juan Carlos Laso; Lunga, Dalton; McCarty, Jessica L.; Pontius, Robert Gilmore; Reinmann, Andrew B.; Rogan, John; Song, Lei; Stoyanova, Hristiana; Ye, Su; Yi, Zhuang Fang; and Estes, Lyndon, "Accounting for training data error in machine learning applied to earth observations" (2020). *Geography*. 56.

[https://commons.clarku.edu/faculty\\_geography/56](https://commons.clarku.edu/faculty_geography/56)

This Article is brought to you for free and open access by the Faculty Works by Department and/or School at Clark Digital Commons. It has been accepted for inclusion in Geography by an authorized administrator of Clark Digital Commons. For more information, please contact [larobinson@clarku.edu](mailto:larobinson@clarku.edu), [cstebbins@clarku.edu](mailto:cstebbins@clarku.edu).







---

**Authors**

Arthur Elmes, Hamed Alemohammad, Ryan Avery, Kelly Caylor, J. Ronald Eastman, Lewis Fishgold, Mark A. Friedl, Meha Jain, Divyani Kohli, Juan Carlos Laso Bayas, Dalton Lunga, Jessica L. McCarty, Robert Gilmore Pontius, Andrew B. Reinmann, John Rogan, Lei Song, Hristiana Stoyanova, Su Ye, Zhuang Fang Yi, and Lyndon Estes

Review

# Accounting for Training Data Error in Machine Learning Applied to Earth Observations

Arthur Elmes <sup>1,2,\*</sup> , Hamed Alemohammad <sup>3</sup> , Ryan Avery <sup>4</sup>, Kelly Caylor <sup>4,5</sup>, J. Ronald Eastman <sup>1</sup>, Lewis Fishgold <sup>6</sup>, Mark A. Friedl <sup>7</sup>, Meha Jain <sup>8</sup>, Divyani Kohli <sup>9</sup>, Juan Carlos Laso Bayas <sup>10</sup> , Dalton Lunga <sup>11</sup>, Jessica L. McCarty <sup>12</sup> , Robert Gilmore Pontius Jr. <sup>1</sup> , Andrew B. Reinmann <sup>13,14</sup>, John Rogan <sup>1</sup>, Lei Song <sup>1</sup>, Hristiana Stoyanova <sup>13,14</sup>, Su Ye <sup>1</sup>, Zhuang-Fang Yi <sup>15</sup> and Lyndon Estes <sup>1</sup> 

<sup>1</sup> Graduate School of Geography, Clark University, Worcester, MA 01610, USA; reastman@clarku.edu (J.R.E.); rpontius@clarku.edu (R.G.P.J.); jrogan@clarku.edu (J.R.); lsong@clarku.edu (L.S.); sye@clarku.edu (S.Y.); lestes@clarku.edu (L.E.)

<sup>2</sup> School for the Environment, University of Massachusetts Boston, Boston, MA 02125, USA

<sup>3</sup> Radiant Earth Foundation, San Francisco, CA, 94105, USA; hamed@radiant.earth

<sup>4</sup> Department of Geography, University of California, Santa Barbara, CA 93013, USA; ravery@ucsb.edu (R.A.); caylor@ucsb.edu (K.C.)

<sup>5</sup> Bren School of Environmental Science and Management, University of California, Santa Barbara, CA 93013, USA

<sup>6</sup> Azavea, Inc., Philadelphia, PA 19123, USA; lfishgold@azavea.com

<sup>7</sup> Department of Earth and Environment, Boston University, Boston, MA 02215; friedl@bu.edu

<sup>8</sup> School for Environment and Sustainability, University of Michigan, Ann Arbor, MI 48109, USA; mehajain@umich.edu

<sup>9</sup> Faculty of Geo-Information Science & Earth Observation (ITC), University of Twente, 7514 AE Enschede, The Netherlands; d.kohli@utwente.nl

<sup>10</sup> Center for Earth Observation and Citizen Science, Ecosystems Services and Management Program, International Institute for Applied Systems Analysis (IIASA), Laxenburg A-2361, Austria; lasobaya@iiasa.ac.at

<sup>11</sup> National Security Emerging Technologies, Oak Ridge National Laboratory, Oak Ridge, TN 37831, USA; lungadd@ornl.gov

<sup>12</sup> Department of Geography and Geospatial Analysis Center, Miami University, Oxford, OH 45056, USA; mccartjl@MiamiOH.edu

<sup>13</sup> Environmental Sciences Initiative, CUNY Advanced Science Research Center, New York, NY 10065, USA; areinmann@gc.cuny.edu (A.B.R.); Hristiana.Stoyanova22@myhunter.cuny.edu (H.S.)

<sup>14</sup> Department of Geography and Environmental Science, Hunter College, New York, NY 10065, USA

<sup>15</sup> Development Seed, Washington, DC 20001, USA; nana@developmentseed.org

\* Correspondence: arthur.elmes@umb.edu; Tel.: +1-304-906-7946

Received: 8 February 2020; Accepted: 18 March 2020; Published: 23 March 2020



**Abstract:** Remote sensing, or Earth Observation (EO), is increasingly used to understand Earth system dynamics and create continuous and categorical maps of biophysical properties and land cover, especially based on recent advances in machine learning (ML). ML models typically require large, spatially explicit training datasets to make accurate predictions. Training data (TD) are typically generated by digitizing polygons on high spatial-resolution imagery, by collecting in situ data, or by using pre-existing datasets. TD are often assumed to accurately represent the truth, but in practice almost always have error, stemming from (1) sample design, and (2) sample collection errors. The latter is particularly relevant for image-interpreted TD, an increasingly commonly used method due to its practicality and the increasing training sample size requirements of modern ML algorithms. TD errors can cause substantial errors in the maps created using ML algorithms, which may impact map use and interpretation. Despite these potential errors and their real-world consequences for map-based decisions, TD error is often not accounted for or reported in EO research. Here we review the current

practices for collecting and handling TD. We identify the sources of TD error, and illustrate their impacts using several case studies representing different EO applications (infrastructure mapping, global surface flux estimates, and agricultural monitoring), and provide guidelines for minimizing and accounting for TD errors. To harmonize terminology, we distinguish TD from three other classes of data that should be used to create and assess ML models: training reference data, used to assess the quality of TD during data generation; validation data, used to iteratively improve models; and map reference data, used only for final accuracy assessment. We focus primarily on TD, but our advice is generally applicable to all four classes, and we ground our review in established best practices for map accuracy assessment literature. EO researchers should start by determining the tolerable levels of map error and appropriate error metrics. Next, TD error should be minimized during sample design by choosing a representative spatio-temporal collection strategy, by using spatially and temporally relevant imagery and ancillary data sources during TD creation, and by selecting a set of legend definitions supported by the data. Furthermore, TD error can be minimized during the collection of individual samples by using consensus-based collection strategies, by directly comparing interpreted training observations against expert-generated training reference data to derive TD error metrics, and by providing image interpreters with thorough application-specific training. We strongly advise that TD error is incorporated in model outputs, either directly in bias and variance estimates or, at a minimum, by documenting the sources and implications of error. TD should be fully documented and made available via an open TD repository, allowing others to replicate and assess its use. To guide researchers in this process, we propose three tiers of TD error accounting standards. Finally, we advise researchers to clearly communicate the magnitude and impacts of TD error on map outputs, with specific consideration given to the likely map audience.

**Keywords:** training data; machine learning; map accuracy; error propagation

---

## 1. Introduction

Recent technological advancements have led to a new era in Earth observation (EO, also known as remote sensing), marked by rapid gains in our ability to map and measure features on the Earth's surface such as land cover and land use (LCLU), e.g., [1,2], vegetation cover and abundance [3], soil moisture [4], infrastructure [5,6], vegetation phenology [7–9], land surface albedo [10–12], and land surface temperature [13,14]. The resulting data are used by an expanding set of disciplines to gain new insights into socioeconomic and environmental dynamics, such as community-level poverty rates [15], changes in surface water [16] and forest cover [17], and carbon accounting [18]. As such, EO is increasingly shaping our understanding of how the world works, and how it is changing.

These breakthroughs are facilitated by several technological advances, particularly the increasing availability of moderate (5–30 m), high-resolution (1–5m, HR), and very high resolution (<1 m, VHR) imagery, as well as new machine-learning (ML) algorithms that frequently require large, high quality training datasets [19–24]. Large training datasets have been necessary for decades in the production of continental and global maps [1,2,25,26]. In the current data-rich era, the impact of training data (TD) quality and quantity on map accuracy is even more relevant, especially for maps generated by data-hungry ML algorithms [27–32]. Errors in these products also impact the veracity of any downstream products into which they are ingested [33]. While progress in algorithmic performance continues apace, standards regarding the collection and use of TD remain uncoordinated across researchers [34]. Additionally, much of the research and development of big data and ML is occurring in industry and the fields of computer science and (non-spatial) data science, leaving a potential knowledge gap for EO scientists [35,36].

The measurement and communication of map accuracy is a mature topic in EO and related fields, with a variety of metrics and approaches tailored to different data types, analyses, and

user groups [37–45]. This includes substantial work to measure error in map reference data (i.e., the independent sample used to assess map accuracy) and account for its impact on map assessment [34,38,46,47]. However, focus on the quality and impacts of TD error has been less systematic. While several efforts have been made to use and evaluate the impact of different aspects of TD quality (noise, sample design, and size) on classifiers [30,32,48–53], much of this work focuses on exploring these issues for specific algorithms [31,48,53,54]. Previous research shows that the impact of TD error can be substantial but varied, suggesting that a more comprehensive approach to this issue is warranted. Furthermore, while TD and map reference data are often collected using the same approaches [55–57] and often subject to the same errors, the existing procedures to minimize and account for map reference errors [34,38,46,47] are not necessarily relevant for quantifying the impacts of TD error. The problems associated with TD error can be summarized as follows:

1. The “big data” era vastly increases the demand for TD.
2. ML-generated map products rely heavily on human-generated TD, which in most cases contain error, particularly when developed through image interpretation.
3. Uncertainty in TD is rarely assessed or reported, and TD are often assumed to have perfect accuracy [30] (which is also common with map reference data [57]).
4. TD errors may propagate to downstream products in surprising and potentially harmful ways (e.g., leading to bad decisions) and can occur without the map producer and/or map user’s knowledge. This problem is particularly relevant in the common case where TD and reference data are collected using the same methods, and/or in cases where map reference data error is not known or accounted for, which is still common [57].

These problems suggest a pressing need to review the issues surrounding TD quality and how it impacts ML-generated maps, and to recommend a set of best practices and standards for minimizing and accounting for those errors, which are the primary aims of this paper. Although map error can also originate from other sources, such as the specific ML classifier selected or the parameterization approach used [31,58,59], we focus solely on issues of input data quality. As such, this paper complements existing work focused on assessing final map accuracy [37–41,44,45].

This paper is organized into four sections. In Section 1, we review current practices in the treatment of TD for categorical and continuous map creation. We also cover map accuracy procedures, given that the two processes are often intertwined and affected by many of the same issues [47], and accuracy assessment procedures are needed to assess the impacts of TD error. In Section 2, we identify the most common sources of TD error and inconsistency. In Section 3, we illustrate the impacts of uncertainty in TD generation with case studies that span a range of typical EO applications: building and road mapping, global surface flux estimates, and mapping agricultural systems. In Section 4, we propose guidelines for (1) best practices in collecting and using TD, (2) minimizing TD errors associated with training sample design error and collection, (3) characterizing and incorporating TD error in final map outputs, and (4) communicating TD error in scientific and public documentation.

### 1.1. Current Trends in Training Data (TD) Collection

A large proportion of remote-sensing projects make some use of TD, typically created either using geolocated in situ data [46,60], by visually interpreting high and/or very high spatial-resolution imagery [26,61,62], or by interpreting the images to be classified/modeled themselves, e.g., [55,56,63,64]. Of these collection methods, HR/VHR image interpretation is increasingly common [65], particularly with the rise in crowdsourcing initiatives [22,66]. As such, mapping is strongly constrained by the creation of TD, which, much like map reference data, are often treated as absolute “truth”, in that their accuracy is assumed to be perfect [30,38,47,67]. However, multiple sources of error are possible and indeed likely in TD, whether collected in situ or via image interpretation [60].

The use of large, data-intensive ML algorithms continues to grow in many fields, including remote sensing. Neural networks (NN) represent an increasingly used class of ML algorithms, with more

complex NNs such as convolutional neural networks (CNN) producing higher output accuracy [68]. While some forms of ML can function effectively with smaller training datasets, the quality of these data is nevertheless critically important [28,31,51]. Additionally, the increasingly popular large-scale, high-complexity NNs require substantially more TD than traditional statistical models, and like many ML approaches are sensitive to noisy and biased data, producing the logistical difficulty of creating very large, “clean” training datasets [69–71].

Partially to address this need, several recent efforts have been devoted to producing extremely large training datasets that can be used across a wide range of mapping applications, and to serve as comprehensive benchmarks [72,73]. Similarly, a recent trend has emerged in large-scale mapping projects to employ large teams of TD interpreters, often within citizen science campaigns that rely on web-based data creation tools [22,74–76].

## 1.2. Characterizing Training Data Error

Due to different disciplinary lineages, terminology associated with the various datasets used to train and evaluate map algorithms is sometimes contradictory or disparate. Here we harmonize terminology by defining four distinct types of data: training, validation, training reference, and map reference. *Training data* (TD) refers to a sample of observations, typically consisting of points or polygons, that relate image pixels and/or objects to semantic labels. *Validation data* are typically a random subset of TD that are withheld and used to fit ML model parameters and internally evaluate performance. *Training reference data* are expert-defined exemplar observations used to assess TD errors during or after data creation. *Map reference data* are independent observations used to assess final map accuracy; while these may be collected using many of the same procedures as the other three datasets [57], they have more stringent design protocols and can only be used to assess the final map product, rather than used iteratively in model or map improvement [57]. Map reference data are often referred to as the test set in ML literature [77], but we use the former term to align with the terminology commonly used by the EO community.

### 1.2.1. Map Accuracy Assessment Procedures

Map accuracy assessment practices and standards are well-established in the EO literature [39,40,45,57,78]. We briefly review these procedures here because they are essential for quantifying how TD error impacts map accuracy. Additionally, the growing use of ML algorithms developed outside of EO has brought with it accuracy assessment practices and terminology that often differ nominally or substantively from those developed for EO, e.g., [57,79,80]. Reviewing EO accuracy assessment standards can, therefore, help to harmonize and improve accuracy assessment practices, while providing necessary context for procedures that can help to account for TD error.

The accuracy of a map is assessed by evaluating the agreement between the values of the mapped variables and those of a map reference variable, and summarizing those discrepancies using an accuracy metric [41,57]. The accuracy metric selected depends on whether the mapped variable is categorical or continuous, since each type of variable has its own foundation for error analysis [81–85]. For categorical variables, this foundation is provided by the confusion matrix, in which rows (but sometimes columns) typically list how many mapped values fall within each category and columns (but sometimes rows) list the distribution of map reference values for each category. In EO, the most widely used metrics calculated from the confusion matrix are user’s accuracy (the complement of commission error), producer’s accuracy (the complement of omission error), and overall accuracy (i.e., the complement of proportion error) [40]. A fuller explanation of accuracy metrics and other aspects of the error matrix can be found in existing publications [37,39,57,81,86–88]. Another widely used measure in EO is the Kappa index of agreement [89], but Kappa varies with class prevalence [90] and inappropriately corrects for chance agreement [57], thus its continued use is strongly discouraged [40,57,91]. There are a number of other categorical accuracy metrics suitable for assessing the accuracy of a binary

categorical variable, such as the F1 score [80], and the true skill statistic [90], which are described in the supplemental materials.

The scatter plot provides the basis for error analysis for continuous variables, wherein deviations between the mapped values plotted on the Y-axis are measured against those of the map reference on the X-axis. Several measures are used to summarize these deviations (see supplementary materials). The root mean squared error (RMSE, also known as root mean square deviation, RMSD) and mean absolute deviation (MAD) summarize deviations along the identity line, also referred to as the 1:1 or  $y = x$  line. RMSE has widespread use, but we recommend caution since it combines MAD with variation among the deviations [92–94]. Another widely used measure is the  $R^2$ , or coefficient of determination, but this measures deviation relative to the linear regression line, rather than the  $y = x$  line [82,92].

Beyond these, there are measures for comparing continuous mapped variables to a binary reference variable, including the receiver operating characteristic (ROC) and the total operating characteristic (TOC) [83,95,96]. The area under this curve (AUC) of an ROC/TOC plot is often used as a single measure of overall accuracy that summarizes numerous thresholds for the continuous variable [96]. There are also metrics for assessing the accuracy of object-based image analysis (OBIA, [97]), which we do not cover here (but see the supplementary information (SI)) because the choice of measure varies according to mapping objectives [65,98].

The creation of the map reference sample is an integral part of the accuracy assessment process and has two major aspects. The first of these is the design of the sample itself (i.e., the placement of sample units), which should be probability-based but can follow several different designs (e.g., simple random, stratified, cluster, systematic) depending on the application and a priori knowledge of the study area [39,57]. The second aspect is the response design, which governs the procedures for assigning values to the map reference samples [39,57]. These include the choice of the sample's spatial and temporal units, the source of the data that the sample extracts from (e.g., high resolution imagery), and the procedure for converting reference data values into map-relevant values [39,57]. For a categorical map in which the reference data source is high-resolution imagery, the map reference sample is assigned labels corresponding to the map legend (e.g., a land-cover scheme) based on a human supervisor's interpretation of the imagery [57].

A key aspect of response design is that map reference data should be substantially more accurate than the map being assessed, even though they are always likely to have some uncertainty [30,39,46,47,57]. This uncertainty should be measured and factored into the accuracy assessment [39,46]. However, in practice this accounting is rarely done, while map reference data uncertainty is also rarely examined [34,38,57]. This tendency is illustrated by Ye et al. [65], who reviewed 209 journal articles focused on object-based image analysis, finding that one third gave incomplete information about the sample design and size of their map reference data, let alone any mention of error within the sample. Errors in map reference data can bias the map accuracy assessment [47,99], as well as estimates derived from the confusion matrix, such as land cover class proportions and their standard errors [46]. To correct for such impacts to map accuracy assessment, one can use published accuracy assessment procedures, including variance estimators, that account for map reference error [38,46,47]. These approaches depend on quantifying errors in the map reference data.

### 1.2.2. Current Approaches for Assessing and Accounting for Training Data Error

Most of the aforementioned considerations regarding map reference data creation largely apply to TD, particularly since map reference data and TD may often be collected together, e.g., [55], provided the former are kept strictly separate to ensure their independence [57]. Considerations regarding TD may diverge with respect to sample design, as TD often need to be collected in ways that deviate from probability-based sampling, in order to satisfy algorithm-specific requirements related to, for example, class balance and representativeness or the size of the training sample [31,51]. Another difference is that map TD can be usable even with substantial error [48,50,51]—although we show in Section 3 that

TD error can propagate substantial map error—whereas map reference data needs to have the highest possible accuracy and its uncertainty should be quantified, as described above [39,46,57].

If the quality of map reference data is often unexamined, TD quality may be even less so. To gain further insight into the level of attention TD receives in EO studies, we reviewed 30 top-ranked research papers published within the previous 10 years that describe land cover mapping studies. (Publications identified by Google Scholar search algorithm results; the search was performed in January 2019, with terms land cover and land use mapping, including permutations of spelling and punctuation. Twenty-seven articles kept after initial screening for relevance—see Table S1 [2,63,64,100–123]). This assessment showed that only three papers explicitly and systematically assessed the quality of the TD used in classification [2,115,122], while 16 made no mention of TD standards at all. Over 75% of these studies used image interpretation, as opposed to in situ data, in either training, accuracy assessment, or both. One-quarter of these papers used unsupervised classifiers in the processing chain to outline training areas, followed by image interpretation to assign labels to the polygons/pixels. Although only a snapshot, this finding suggests that key details regarding the design and collection of TD (and even map reference data) is lacking in the EO literature.

Even though TD quality appears to be largely unreported, efforts have been made to examine how TD error can impact ML-based classifications, typically within the context of evaluating specific algorithms. For example, research examining the effectiveness of random forests [124] for land-cover classification also evaluated their sensitivity to TD error, sample size, and class imbalance [48,51,125]; similar research has been conducted for Support Vector Machines (SVM) [28,32,52]. Several studies comparing multiple ML algorithms also compared how each reacted variations in TD sample size and/or error [50,59,126,127]. Maxwell et al. [31] touch on a number of these TD quality issues in an even broader review of ML algorithms widely used in EO classification but excluding newer deep learning approaches.

Beyond these examples, several studies have focused more explicitly on how to train ML-algorithms for remote sensing classification when TD error is present. Foody et al. [30] conducted tests to examine how two different types of TD labeling error impacted land-cover classifications, with a primary interest in SVM. Similarly, Mellor et al.'s [48] study measured uncertainty introduced by TD error in a random forest classifier, with specific focus on class imbalance and labeling errors. Swan et al. [49] examined how increasing amounts of error introduced into the TD for a deep-learning model impacted its accuracy in identifying building footprints. These studies collectively demonstrate that TD has substantial impact on ML-generated maps. They also reveal that there is no standard, widely accepted practice for assessing TD error, which, similar to map reference data, is generally not reported and thus implicitly treated as error-free [30].

## 2. Sources and Impacts of Training Data Error

In the following two sections we describe the common causes of TD error and explore its potential impacts. To describe these causes, we divide the sources of TD error into two general classes: (1) errors stemming from the design of the training sample, including some aspects of sample and response design that are shared with standards for the collection of map reference data (see 1.2.1 above), and (2) errors made during the collection of the training sample, including additional elements of response design such as the process of digitizing and labeling points or polygons when interpreting imagery or when collecting field measurements. In addressing the impacts of error, we provide a summary of potential problems, and then two concrete case examples for illustrative purposes.

### 2.1. Sources of Training Data Error

#### 2.1.1. Design-Related Errors

With respect to TD sampling design, errors primarily stem from failures to adequately represent the spatial-temporal-spectral domains of the features of interest in the manner most suited to the



specific ML algorithm being used [53]. This problem may be exacerbated in cases where TD are collected exclusively using the same rigorous probability-based specifications used to collect map reference data, which may be overly restrictive for the purposes of TD collection. While the use of such standards to collect TD may be possible provided that there is a large enough data set (e.g., a large benchmark data set), smaller training data sets and/or cases of geographically sparse target classes/objects will benefit strongly from the increased flexibility afforded to TD collection standards, which are less restrictive than those for map reference data (e.g., allowing for purposive rather than purely probabilistic sampling). A lack of geographic representation of the phenomena of interest results in a disparity between the distribution of TD compared to the true distribution of the mapped phenomenon in geographic and/or feature space [28–31]. This problem is highly relevant in ML approaches, which are sensitive to TD quality, including class balance, labeling accuracy, and class comprehensiveness relative to the study area’s true composition [30].

Temporal unrepresentativeness is also a common source of error in the response design of TD, due to the prevalence of image interpretation as a source for TD. In this case, error arises when obsolete imagery is interpreted to collect training points or polygons and their associated labels [39,61]. The problem is illustrated in Figure 1, which contrasts smallholder fields that are clearly visible in a satellite base map (Bing Maps) with ground data collected in 2018. Center pivot fields were installed after the base map imagery was collected, but before ground data collection, causing a temporal mismatch between the base map and the in situ data. Labels generated from the base map would therefore introduce substantial error into an ML algorithm classifying more recent imagery. New HR/VHR satellites that have more frequent acquisitions (e.g., PlanetScope [128]) can help minimize such temporal gaps for projects that are designed to map present-day conditions (e.g., 2018 land cover), but cannot solve this problem for mapping projects covering earlier time periods (i.e., before 2016). The same can be said for aerial and unmanned aerial vehicle acquisitions, which are typically limited in geographic and temporal extent [129]. While hardcopy historical maps can help supplement temporal data gaps, these data sources come with their own problems, such as errors introduced during scanning and co-registration, and unknown production standards and undocumented mapping uncertainties.



**Figure 1.** An example of potential training data error that can arise when image interpretation is conducted on older imagery. The underlying imagery is from Bing Maps, which shows smallholder agricultural fields near Kulpawn, Ghana. The white polygons were collected by a team of mappers (hired by Meridia) on the ground using a hand-held Global Positioning System (GPS) in 2018. The smallholder fields were replaced by larger center-pivot irrigation fields sometime after the imagery in the base map was collected.

Spatial co-registration can be a substantial source of response design error when training with HR and VHR commercial satellite imagery. Due to their narrow swath widths, HR/VHR sensors are often

tasked, resulting in substantially off-nadir image acquisitions [61]. Due to large view zenith angles and the lack of adequate digital elevation models, side overlapping imagery for stereo photogrammetry, or other relevant control points, HR/VHR imagery often does not meet the same orthorectification standards as coarser resolution, government-operated satellites [130–132]. When integrating HR/VHR imagery acquired at different azimuth and elevation angles, features such as building roofs show offsets similar to those caused by topography. These offsets are particularly problematic for (a) training repeated mappings of the same features, and/or (b) when using an existing vector dataset such as OpenStreetMap (OSM) as TD [133–135].

TD collected by interpreting HR/VHR imagery is often co-registered with the coarser resolution imagery used as ML model data. This creates a potential spatial resolution conflict because the minimum mapping unit (MMU), i.e., the relationship between image objects and pixel size, may be different in the two imagery data sets. This potentially leads to situations in which objects delineated as spectrally homogenous areas in HR/VHR imagery are part of mixed pixels in moderate- or coarse-resolution model imagery. This mismatch is similar to the concept of H-resolution versus L-resolution scene models proposed by Strahler et al. [136]; in H-resolution models, the objects of interest are substantially larger than the pixel size, and vice versa for L-resolution models. The incorporation of mixed pixels may degrade classification model performance, or at least introduce undesired spectral variability within classes [127,137,138]. This situation may be alleviated by displaying both HR/VHR imagery and/or other ancillary datasets as well as coarser model imagery during training data creation [139,140]. However, such practices may not be possible when training data are taken from previous research projects, or when they are to be applied in the context of time series analysis, in which spatial features change over time, e.g., [141].

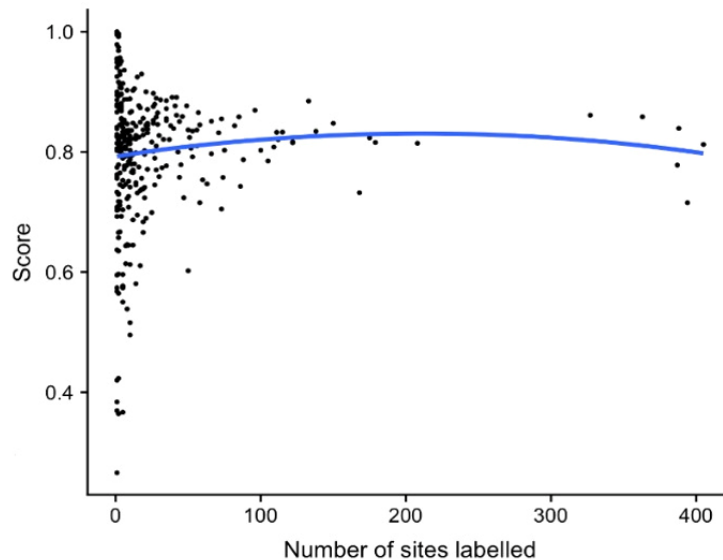
Similar spatial resolution and scaling issues must be dealt with when combining in situ measurements with satellite observations for continuous variables. Field-collected data often cannot practically cover the entire area of a pixel in the model data, especially for moderate or coarse-resolution imagery, and can thus induce scaling errors related to the modifiable areal unit problem [142,143]. Spatial representativeness assessments and interpolation methods are used to limit this problem for operational EO science products [144–147], but this issue is likely to be a source of error for most in situ TD samples.

Another design-related problem arises from large-scale data collection initiatives that are becoming increasingly common due to the expanding extent of modern EO analyses, e.g., [148]. These efforts, often conducted via crowdsourcing campaigns, typically enlist citizens to collect data via a web-based platform, e.g., [66,149–151]. Examples include OSM, Geo-Wiki [66], Collect Earth [152], DIYLandcover [150], and FotoQuest Go [153]. In cases where the resulting data might be purely voluntary [76], the resulting sample may lack spatial representativeness due to uneven geographic contributions [28,154].

### 2.1.2. Collection-Related Errors

There are several common forms of error that occur when collecting both TD and map reference data. The first of these are errors of interpretation [39], which are mistakes created in the process of manual image interpretation. Image interpretation is widely used to generate TD, but often this technique leads to inconsistent labels between interpreters for the same areas of interest [34,37,99,155]. Interpreters may lack experience in the task or be unfamiliar with the context of the study area, e.g., [156]. In an unusually thorough analysis of error in image interpretation, Powell et al. [99] showed that inter-interpreter agreement was on average 86% but ranged from 46 to 92%, depending on land cover. This research, which relied on trained image interpreters, concluded that transitional land cover classes produce substantial interpretation uncertainty, which is particularly problematic since much land cover mapping effort is directed towards change detection. Another image interpretation study that used a crowdsourcing platform found that interpreters' average accuracy in digitizing crop field boundaries in high-resolution imagery was ~80%, based on comparisons against training

reference data [150]. This result held true whether the interpreters mapped several hundred sites or <50 (Figure 2), indicating that increased interpreter experience does not necessarily eliminate labeling error, even when analysts are highly seasoned [99]. These findings underscore the need to assess uncertainty in TD, as well as map reference data, using predefined training reference data or inter-interpreter comparisons [46,60,99,157,158].

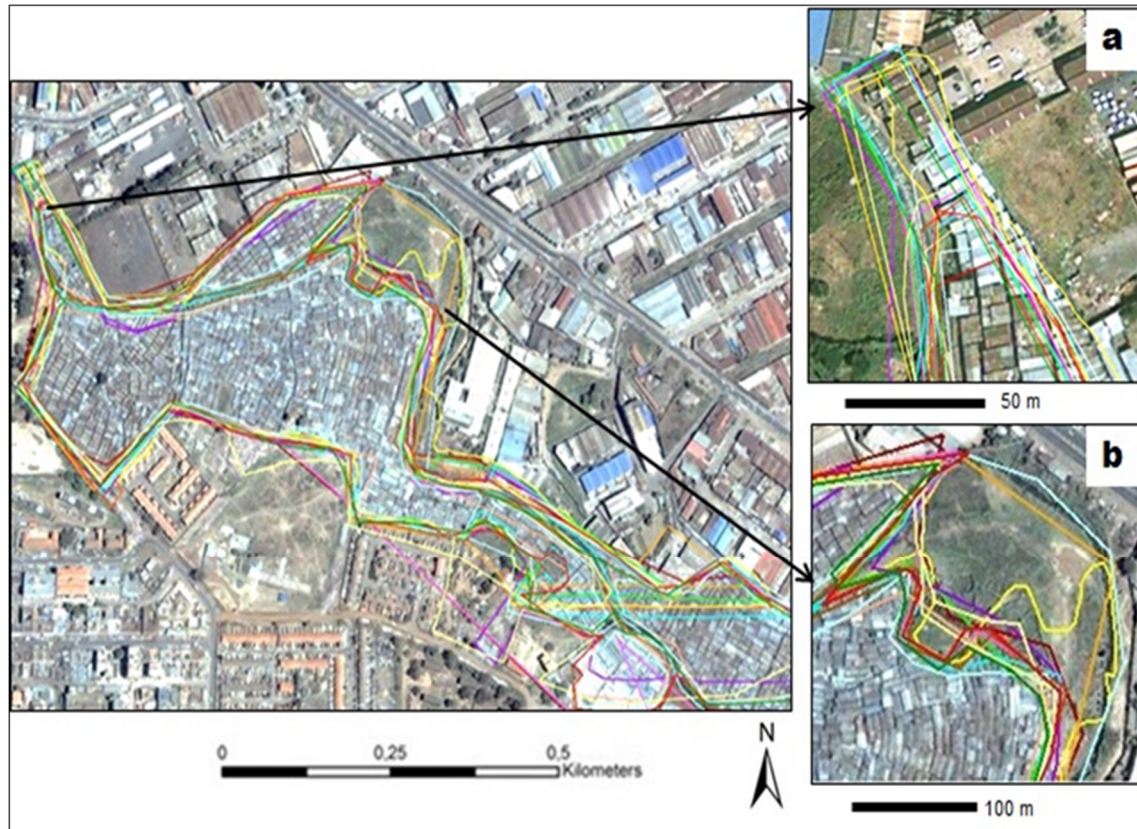


**Figure 2.** Number of sites mapped per worker versus the average score received at reference sites, where workers' maps were compared to reference maps using a built-in accuracy assessment protocol within a crowdsourcing platform for collect cropland data [150].

Labeling error may also result from inadequate or poorly communicated semantic class definitions [159,160], particularly when identifying human land use as opposed to biophysical land cover [161]. This is especially evident in urban environments, which exhibit high spatial and spectral heterogeneity (even within HR/VHR imagery [162]), and are also semantically vague (i.e., hard to define) even at the ground level. For example, Figure 3 shows a typical example of TD collection for mapping informal settlements (i.e., slums), in Nairobi, Kenya, in which several trained interpreters separately delineate the same area [163]. Because slums may be defined by sociodemographic factors in addition to spatial and spectral properties, TD creation for such areas is prone to error stemming from semantic issues [160]. Complex classes such as slums may exhibit high variability between study areas, as local idiosyncrasies link the definition of slums to different physical, remotely observable characteristics. These characteristics make it hard to develop a generalizable mapping capability for land uses such as informal settlements. These results further illustrate the importance of consensus mapping for image interpretation, particularly for spatially, spectrally, or temporally heterogeneous LCLU classes, which may have vague or regionally idiosyncratic semantic definitions.

Categorical mapping projects typically define a crisp set of non-overlapping categories, rather than a fuzzy set [164,165]. However, many human and natural land covers exhibit continuous gradation between classes, implying that crisp map legends will necessarily cause semantic ambiguity when image pixels in transitional areas are labeled [166,167]. This problem is particularly acute with moderate- and coarse-resolution imagery [26]. Local variance is highest when scene objects approximate the spatial dimension of the image resolution, leading to poor classification accuracy [168]. While substantial research has been devoted to the issue of mixed pixels [85,137,138,169–171], crisp categories are still often relied on during the training and testing phases of image classification [172]. Alternative approaches based on fuzzy set theory are available, but have seen limited adoption [165,173]. Labeling errors can also arise if analysts are not properly trained regarding class definitions, or by the failure to capture comprehensive metadata while collecting TD in the field or during image interpretation. Lack

of TD metadata is particularly problematic in the context of difficult-to-determine labeling cases, or when there is potential confusion between spectrally, spatially, or semantically/conceptually similar classes [161]. Such inadequacies limit the analysis of TD error and, therefore, the ability to account for error propagation.



**Figure 3.** The challenges of mapping slum extent from image interpretation in Nairobi, Kenya. Each colored line indicates a different analyst's delineation of the same slum, illustrating semantic confusion. Adapted with permission from Kohli et al. [163].

Collection-related errors may be particularly acute in large-scale crowdsourcing campaigns or citizen science initiatives, which are increasingly valued for mapping projects due to their larger size and cheaper acquisition costs [22,66,150,151]. Such datasets are often collected rapidly and entail labeling many observations over a short period of time by participants who are not domain experts [153,174]. In such cases, label quality is a function of interpreter skill, experience, contextual knowledge, personal interest, and motivation for involvement in the data collection [22]. Errors can be exacerbated if interpreters are inadequately trained or unfamiliar with the study area, or lack experience with EO data and methods. For example, delineation of different classes of urban land use may be extremely difficult without the benefit of local knowledge [160]. Furthermore, image interpretation is complicated when participants are required to interpret HR/VHR satellite imagery collected over multiple sensors, on different acquisition dates, with varying quality (e.g., cloud cover percentage and atmospheric correction), and/or with varying view/sun angles [175]. Inadequate or confusing user interfaces may also lead to error [22,160]. Once crowdsourced/citizen science data have been post-processed for noise, they can be highly detailed and spatially extensive [66,69–71]. Nevertheless, quality problems in such datasets can be particularly hard to find and clean and are thus an important source of TD error that may propagate through ML algorithms into map outputs [57,151,176]. Therefore, these data should be used more cautiously than expert-derived TD.

Errors also arise in in situ TD, caused by measurement error, geolocation inaccuracy, or incorrect identification of relevant objects (e.g., vegetation species), for example [177]. In addition to these factors, some feature types may also be difficult to discern on the ground [30]. Aside from these problems, there are many sources of technologically induced errors, such as defects in the software or hardware of measurement devices, user input error, or calibration errors (e.g., in spectro-radiometers or other equipment). However, accounting for quantitative measurement error is more straightforward than thematic TD creation. Textbook tools to quantify measurement error are widely available, and in situ data collection procedures often include inter-analyst measurement comparison [178,179].

## 2.2. Impacts of Training Data Error

TD errors carry through to impact the map production process and outcomes. From a design perspective, the size and class composition of TD is particularly impactful on ML algorithms, which are susceptible to overfitting and class imbalance problems [31,73]. Additionally, the assumption of representativeness of training pixels is often overstated, and many TD may in fact not be generalizable to broader scales (discussed by Tuia et al. [154]). TD errors arising from the collection process also impact map quality. Both design- and collection-related errors may be particularly hard to discern, or quantify in absolute terms, if the error in the map reference data errors are unknown.

Several studies reviewed in Section 1.2.2 provide insight into how much TD error can impact ML-generated land-cover maps, focusing on aspects of sample size and balance (design-related errors) and labeling error (collection-related error). This work shows that the impact of each error source varies according to the algorithm used. For example, SVMs were relatively insensitive to changes in sample size, with accuracy dropping by only 3%–6% under TD size reductions of 85–94% [28,180]. Random forests (RF) also proved robust to TD sample size, showing slightly higher accuracy drops of ~4–10+% when TD was reduced by 70–99% [48,51,180]. Sample size also impacts the certainty of RF classification by lowering the mean margin (a measure of certainty related to the number of class votes) by ~50% for sample size reductions of 95% [48]. In contrast to SVM and RF, maps classified with single decision trees are highly affected by TD size, with 13% accuracy loss for TD reductions of 85% [28], and up to 50–85% loss with TD size reductions of 50–70% [51,59]. NNs show varying responses to sample size, depending on their algorithmic design: one NN based on adaptive resonance theory showed accuracy reductions of ~30% to ~65% when TD samples were halved [59], while a feed-forward NN lost just 2% accuracy when TD was reduced by 85% [28].

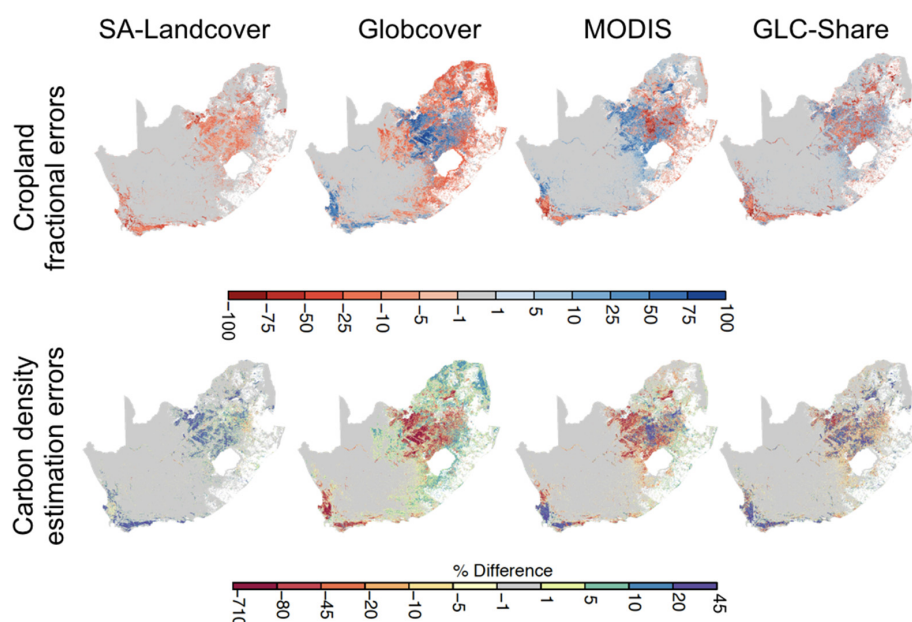
Classifiers are also sensitive to class balance within the training data. For example, the accuracy of RF-generated maps declined by ~12% to ~23% and classification confidence fell ~25% to ~50% when TD class balances were highly skewed [48]. Notably, the ranges in these accuracy and confidence declines were attributable to differing TD sample sizes, showing the synergistic effect of sample size and class balance sensitivities. Maxwell et al. [31] provide a more comprehensive review of class imbalance for RF, SVM, NN, and k-nearest neighbors (kNN) classifiers, finding that all models were sensitive to class imbalance, but the accuracy impact was largest for rare classes, as opposed to overall map accuracy.

The impact of TD labeling errors, also referred to as noise, varies substantially between mapping algorithms. SVMs and closely related derivatives appear least sensitive to mislabeling. SVMs lost just 0–5% in land-cover classification accuracy when 20–30% of TD samples were mislabeled either randomly or uniformly across classes [30,52,126]. Relative vector machines (RVMs) were even less sensitive under these conditions (2.5% accuracy loss for 20% mislabeling [30]), and an SVM designed specifically for handling noisy TD (context-sensitive semi-supervised SVM) was even more robust (2.4% reduction in kappa for 28% mislabeling [52]). However, the impact of TD noise was greater for all three models when mislabeling was confined to specific classes. SVMs lost 9% accuracy and 31% kappa when 20–28% of samples in spectrally similar classes were mislabeled [30,52]. The RVM showed a 6% accuracy loss [30], and the specialized SVM showed a 12% kappa reduction [52] under the same conditions. As with sample size, RF is the next least sensitive to TD noise [48,51]. Mislabeling 25% of

TD samples reduced RF accuracy by 3–7% for a binary classifier and 7–10% for a multiclass model, with the ranges in accuracy loss also varying according to TD sample size [48]. Classification certainty was more heavily impacted by label error, dropping by 45–55%, as measured by the mean margin [48]. Other classification models showed larger impacts due to label noise, including 11–41% kappa declines for a kNN (28% label noise [52]), and 24% [126,181] and 40–43% accuracy loss for a kernel perceptron and NN, respectively, when each is trained with 30% of TD labeled incorrectly [59,126,181]. Single decision-tree models were most sensitive to label error, registering 39% to nearly 70% accuracy declines for 30% label noise [59,126,181].

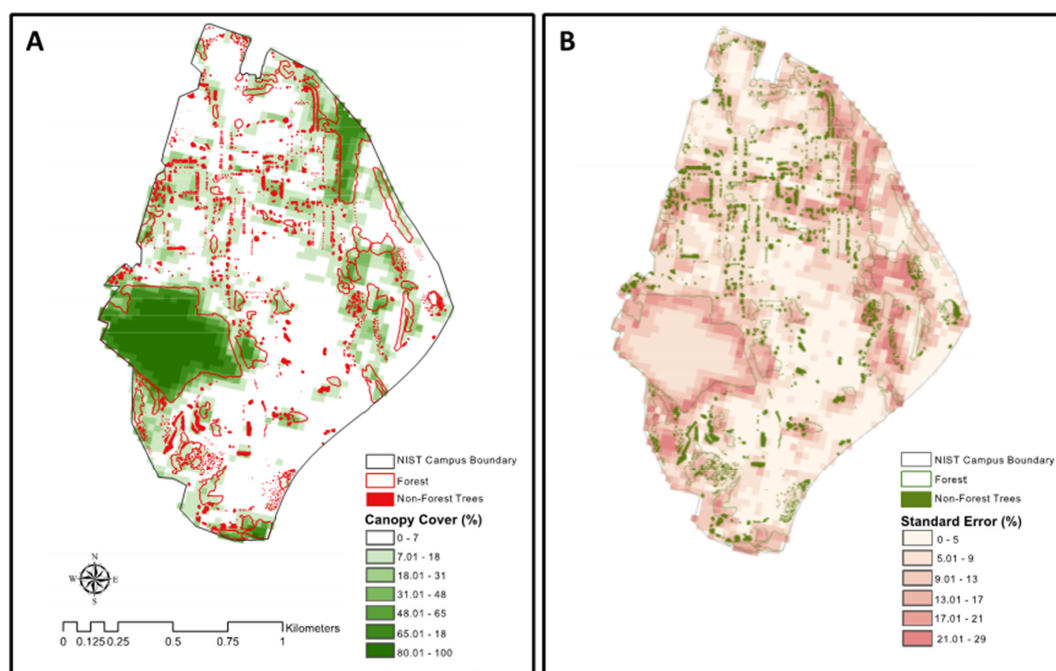
The research described above provides substantial information on how TD error can impact the accuracy and certainty of older-generation ML classifiers. Further understanding of the consequences of these errors can be inferred from literature examining the impact of errors in map reference data. Map reference errors can substantially bias areal estimates of land-cover classes, as well as the estimation of variance in those classes, particularly when examining land-cover change [46,182,183]. While methods exist to incorporate map reference data error into map accuracy assessments and area estimates [38,46,47], and also to account for TD uncertainty in assessing classifier accuracy [48], there has been little work that shows how to address both TD and map reference error.

Less information is available regarding the ways in which TD error may propagate beyond the map it initially creates. Initial research by Estes et al. [33] examined how error propagates from a primary land-cover map into subsequent derived products. This work used a high-quality reference cropland map to quantify the errors in 1 km cropland fractions derived from existing land cover datasets and measured how these errors propagated in several map-based analyses drawing on cropland fractions for inputs. The results suggest that downstream errors were in some instances several fold larger than those in the input cropland maps (e.g., carbon stock estimates, Figure 4), whereas in other cases (e.g., evapotranspiration estimates) errors were muted. In either case, the degree to which the error magnifies or reduces in subsequent maps is difficult to anticipate, and the high likelihood that error could increase means that any conclusions based on such land cover-derived maps must be treated with caution when error propagation is not quantified. This analysis suggests how TD errors might impact the maps they generate and provides a potential method for quantifying their impacts on map accuracy.



**Figure 4.** An examination of how error in pixel-wise cropland fractional estimates (expressed as a percentage, top row) can propagate error (expressed as a percentage) in maps that use land-cover data as inputs, such as estimates of carbon density (bottom row). Figure adapted from Estes et al. [33].

The impact of map input errors can also be seen in the practice of using well-known standard datasets, such as the National Land Cover Map (NLCD, [184]), to map quantities of interest, such as urban tree canopy biomass. Urban trees play a crucial role but in regional carbon cycles [185–187] but are often omitted from EO studies of carbon dynamics, e.g., MODIS Net Primary Productivity [188]. As urban lands are expected to triple between 2000 and 2030 [189,190], the need to factor them into carbon accounting is pressing, but remotely mapping urban tree cover is limited by (a) spatial resolutions that are too coarse for highly variable urban landscapes and (b) TD that are often biased to forested, agricultural, and other rural landscapes. For these reasons, the Landsat-derived NLCD Percent Tree Cover (PTC) product [191], which estimates canopy cover at 30-m resolution across the US, provides a practical input for empirical models to map tree biomass. However, previous studies have shown that this product shows higher uncertainty in urban areas [191] and has a tendency to underestimate urban canopy cover compared to high resolution datasets. Therefore, to quantify the potential impact of NLCD PTC error on canopy biomass estimates, we compared the accuracy of the NLCD PTC dataset to canopy cover estimates derived from manually digitized VHR Imagery for a suburb of Washington, D.C., USA. We found that NLCD PTC underestimated canopy cover by 15.9%, particularly along forest edges (Figure 5) where it underestimated canopy cover by 27%. This discrepancy is particularly important in heterogeneous urban landscapes, where forest edges comprise a high proportion of total forest area. Scaling field data from forest plots to the entire study yielded an estimate of 8164 Mg C stored in aboveground forest biomass, based on our manually digitized canopy cover map, compared to only 5960 Mg C based on the NLCD PTC. This finding indicates the significance of these map errors for carbon accounting, as temperate forest carbon storage and rates of sequestration are much larger (64% and 89%, respectively) than in forest interiors [192]. Quantifying errors in the NLCD is thus important for correcting subsequent estimates trained on these data.



**Figure 5.** Spatial variations in canopy cover (A) and uncertainty in canopy cover estimates (B) in forested and non-forested areas of the heterogeneous suburban landscape of the National Institute of Standards and Technology campus in Gaithersburg, Maryland. Percent canopy cover at a 30-m resolution from the commonly used National Land Cover Database (NLCD) Percent Canopy Cover product (and its uncertainty) is superimposed over a high-resolution map of forested areas (hollow outlined polygons) and non-forest trees (e.g., street trees; solid polygons) that were manually mapped using <1-m resolution Wayback World Imagery. Note the lower estimates of percent canopy cover along forest edges (A) and the associated higher levels of uncertainty (B) using the NLCD product.

These brief examples help illustrate the potential problems of TD error, but the range of potential impacts is as varied as the number of mapping projects underway across academic research, commercial operations, and the public sphere. To represent the growing set of remote-sensing applications in which TD error may be encountered, we present a set of case studies below. To help lay a common framework, we show a typical methods sequence for a ML-based remote-sensing analysis in Figure 6, which also helps clarify the terminology used in this paper. The figure shows the various sources and implications of error in the modeling and mapping process, beginning with issues in the data sources and sample design, and continuing through-model training, validation, and ultimately in map accuracy assessment.

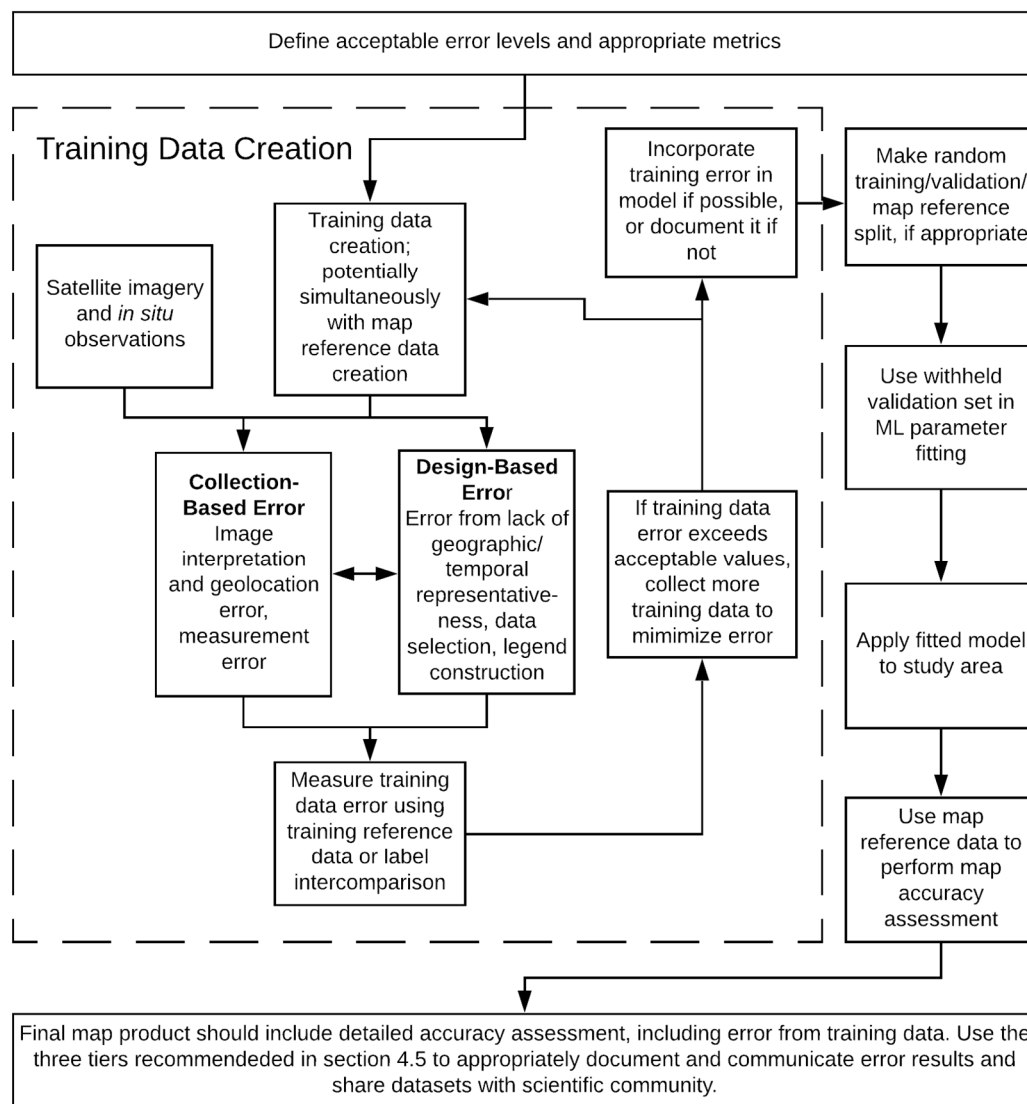


Figure 6. Flow chart of typical workflow for machine-learning applications in Earth observation data.

### 3. Case Studies

To better illustrate the potential impact of TD error, we provide several case studies across different mapping applications that represent the broad range of ML-based mapping and modeling applications that rely on TD.



### 3.1. Infrastructure Mapping

#### 3.1.1. Incorporating Noisy Training Label Data

Automated building footprint detection is an important but difficult mapping task, with potential benefits for a wide range of applications. The following case study illustrates the use of Raster Vision (<https://rastervision.io/>), an open source deep learning framework, to train several models for automated building detection from high resolution imagery (Additional detail available at: <https://www.azavea.com/blog/2019/08/05/noisy-labels-deep-learning/>). These models perform best when trained on a large number of correctly labeled examples, usually generated by a paid team of expert labelers. An alternative, less costly approach was conducted in which a building segmentation model was trained using labels extracted from OSM. However, the labeled training polygons generated from OSM contain errors: some buildings are missing, and others are poorly aligned with the imagery or have missing details. This provides a good test case for experimentation on how noise in the labels affects the accuracy of the resulting model.

To measure the relationship between label noise and model accuracy, the amount of label noise was varied while holding all other variables constant. To do this, an off-the-shelf dataset (the SpaceNet Vegas buildings data set) was used in place of OSM, into which label errors were systematically introduced. Missing and imprecisely drawn building errors were systematically introduced to this relatively large training data set (~30,000 labeled buildings) (<https://spacenetchallenge.github.io/datasets/spacenetBuildings-V2summary.html>), and then the resulting model accuracy was measured. The experimental design consisted of two series of six datasets each, with random deletion or shift of buildings at increasing probabilities and magnitudes, respectively. For each dataset, a UNet semantic segmentation model with a ResNet18 backbone was trained using the fastai/PyTorch plugin for Raster Vision (<https://github.com/azavea/raster-vision-fastai-plugin>). These experiments, including data preparation and visualization, can be replicated using code at [https://github.com/azavea/raster-vision-experiments/tree/master/noisy\\_buildings\\_semseg](https://github.com/azavea/raster-vision-experiments/tree/master/noisy_buildings_semseg).

Figure 7 shows the ground truth and predictions for a variety of scenes and noise levels, showing that the quality of the predictions decreases with the noise level. The background and central portions of buildings tend to be predicted correctly, whereas the outer periphery of buildings presented a greater challenge. These results are quantified in Figure 8, which shows F1, precision, and recall values for each of the noise levels below (see Table S2 for terminology description). The precision falls more slowly than recall (and even increases for noisy drops), which is consistent with the pattern of errors observed in the prediction plots. Pixels that are predicted as building tend to be in the central portion of buildings, leading to high precision.

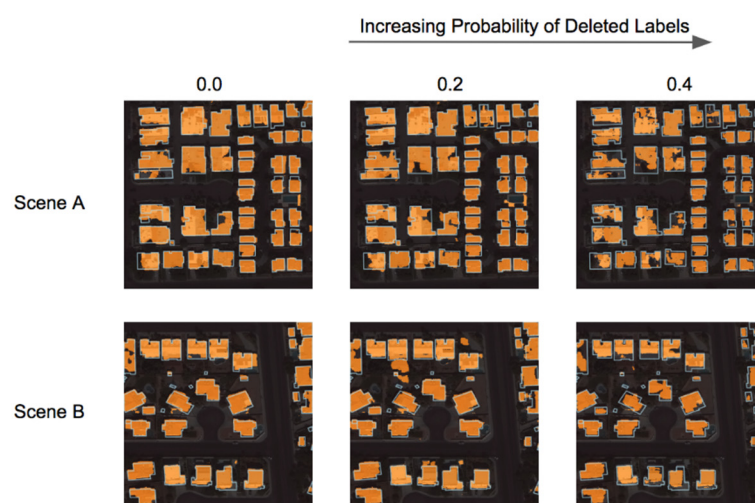
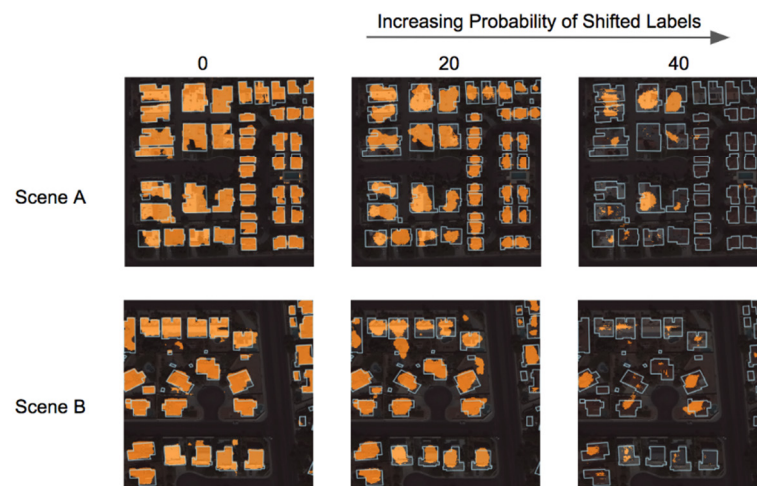
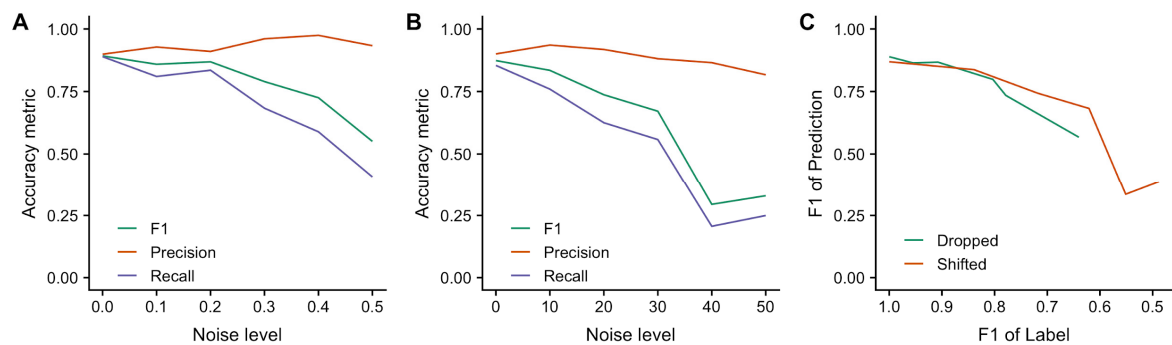


Figure 7. Cont.



**Figure 7.** Predictions of the model trained on different noisy datasets. Each row shows a single scene over different noise levels. The top two rows show noisy drops, while the bottom two rows show noisy shifts. The ground truth is outlined in light blue, and the predictions are filled in orange.



**Figure 8.** The precision, recall, and F1 scores across different noise levels are shown for the cases in which labels are randomly dropped (A) or randomly shifted (B). Panel (C) compares how prediction quality changes as noise increases for dropped and shifted labels, measured by F1 of the labels and prediction.

In panels (A) and (B) of Figure 8, the x-axis shows the noise from randomly dropped and randomly shifted labels, respectively. Panel (C) combines the effects of noisy deletions and noisy shifts on accuracy in a single graph, showing F1 of the labels on the x-axis and F1 of the prediction on the y-axis. The F1 score of the noisy versus ground truth labels is a function of the pixel-wise errors; this metric has the benefit of measuring the effect of noise on error in a way that is comparable across datasets and object classes. For instance, a noisy shift of 10 in a dataset with large buildings might result in a different proportion of erroneous label pixels than in another dataset with small buildings. From this, panel (C) shows that while some of the shifted datasets have a greater level of noise, the prediction F1 scores are similar between the two series when the noise level is similar.

These results present a small step toward determining how much accuracy is sacrificed by using TD from OSM. Preliminary results indicate that accuracy decreases as noise increases, and that the model becomes more conservative as the noise level increases, only predicting central portions of buildings. Furthermore, the noisy shift experiments suggest that the relationship between noise level and accuracy is non-linear. Future work will quantify the functional form of this relationship, and how it varies with the size of the training set. Some preliminary work toward this goal has been described in Rolnick et al. [193], which focuses on image classification of Imagenet-style images.

One limitation of these results is that the magnitude of error in OSM for most areas is unknown, making it difficult to predict the effect of using OSM labels to train models in a generalized, global sense. Noisy error in OSM can be estimated by measuring the disparity between OSM labels to clean labels, such as the SpaceNet labels used here, providing a local estimate of OSM noise. A more general

but less rigorous approach is to roughly estimate the noise level by visually inspecting the labels in OSM, and comparing to Figure 7, which shows examples of the labels at different noise levels.

### 3.1.2. Detecting Roads from Satellite Imagery

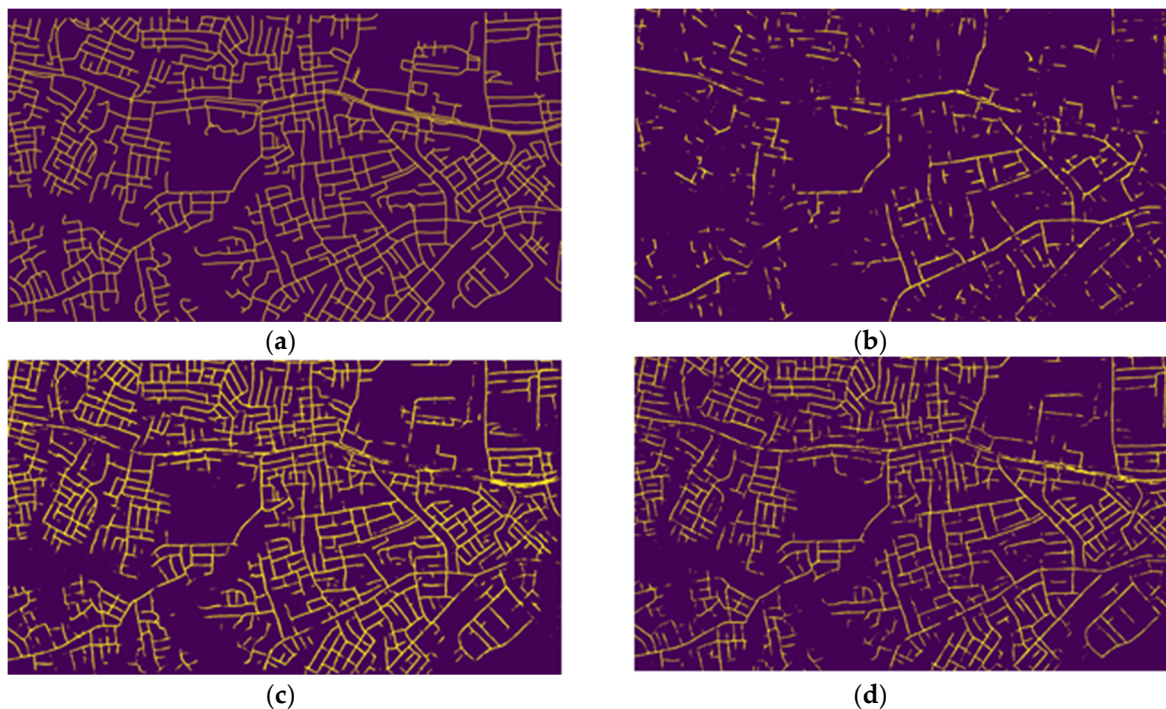
Road networks constitute a critical geographical data layer used to assist national decision makers in resource allocation, infrastructure planning, vaccination campaigns, and disaster response, among others. However, accurate and up-to-date road networks are not available in many developing countries. High-resolution satellite imagery, paired with deep-learning methods, provides the capacity to detect and map roads at large spatial scales. This important goal, however, is dependent on availability of local high-quality TD.

To evaluate the impact of local TD availability on predicted road network accuracy, a study was carried out in Kumasi, Ghana [194]. Two datasets were used to train ML models: (1) the SpaceNet (<https://spacenetchallenge.github.io/>) dataset [195] in Khartoum, Sudan, and Las Vegas, USA, and (2) OSM data in Kumasi, Ghana. The SpaceNet Dataset includes high quality road labels with human expert validation, but unfortunately was not available in Kumasi, Ghana. Therefore, the latter study site relied on OSM data, consisting of crowdsourced labels with no accuracy assessment or expert validation. A series of experiments were carried out to assess the feasibility of using transfer learning, using the Raster Vision Python library for training and evaluation. For all MobileNet V2 models introduced in the following list, the image chip size was set to  $300 \times 300$  pixels, and the training/validation split was 80/20.

The Las Vegas Model was trained and validated on SpaceNet data in Las Vegas and produced very high accuracy predictions. However, when this model was used in Kumasi, it predicted very few roads, with only scattered road segments. The Khartoum model was also trained using SpaceNet data in Khartoum. The Kumasi model used Maxar WorldView-3 imagery and labels from OSM as input. OSM was used to test the quality of crowdsourced labels in training a road detection model. The Khartoum Model was then fine-tuned on OSM labels in Kumasi for three different steps of 100 K, 50 K and 10 K. All models used the same hyperparameters, to isolate the role of TD on model performances.

To validate the models' performance using an independent dataset, a set of expert labels was generated over a small part of Kumasi. Figure 9 shows the region with human expert data vetting, along with the three model predictions. The Las Vegas model is excluded from this figure as it does not have any meaningful prediction in Kumasi. Quantitative performance metrics were calculated using the expert-created labels, to which the models had been blind during training. The results indicate that, as shown by Figure 9, the F1 score for roads was substantially higher for the Kumasi model (0.6458) than when using the Khartoum model (0.3780). However, by retraining and fine-tuning the Khartoum model, the F1 score for roads increased to 0.6135. The full accuracy results for this experiment are presented in Table S3, and prediction maps are shown in Figure S1.

Based on these results, it is concluded that: (1) lack of diverse TD significantly limits the geographic applicability of models, as the types, surfaces, and arrangements of roads varies substantially between regions; (2) regional training datasets are essential for the model to learn the feature of roads in that region; and (3) transfer learning from a reasonably similar geography can help train models.



**Figure 9.** (a) Labels generated by experts for validation. (b) Predictions from the Khartoum model. (c) Predictions from Kumasi model. (d) Predictions from Khartoum model retrained in Kumasi with 10 K steps.

### 3.2. Global Surface Flux Estimates

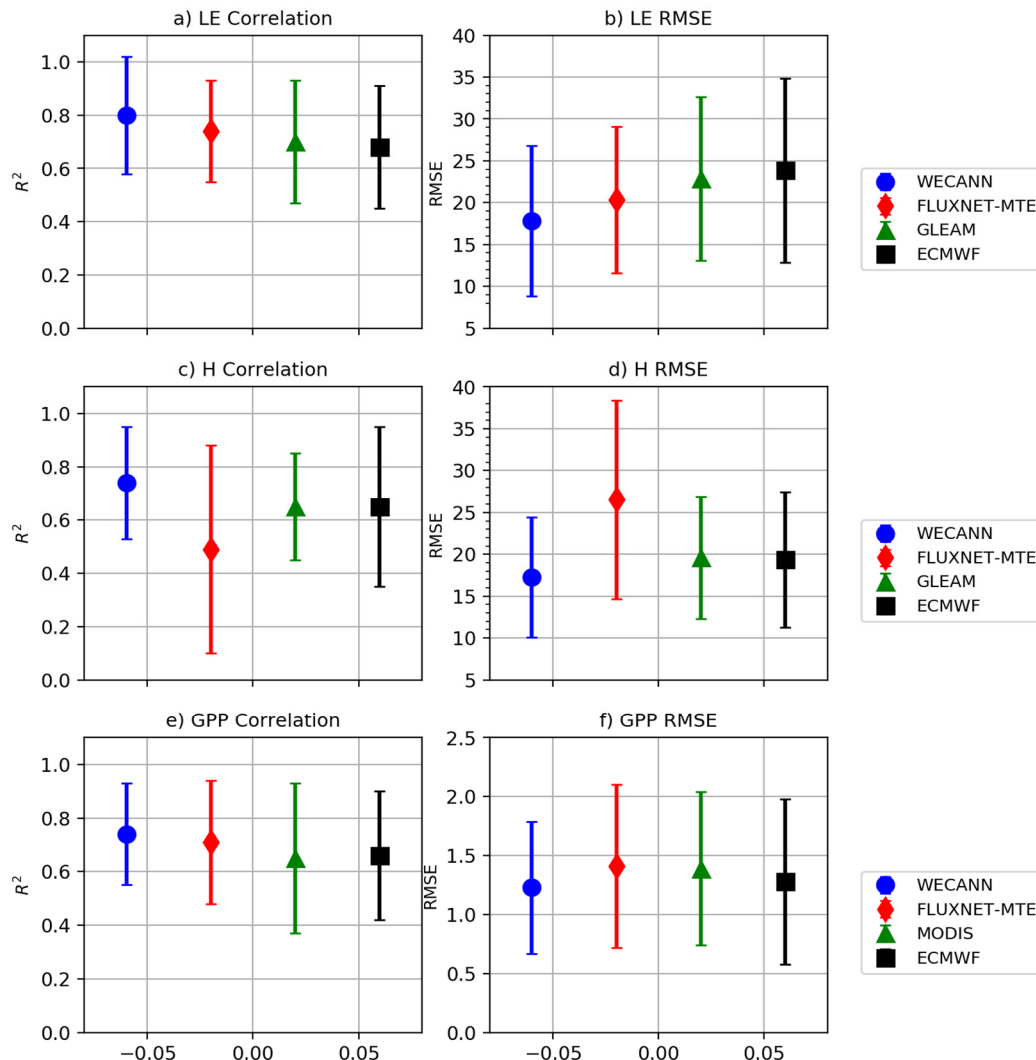
Fluxes at the land–atmosphere boundary play a key role in regulating water, carbon and energy cycles. These fluxes include latent heat flux (LE), sensible heat flux (H), and gross primary production (GPP). While these fluxes cannot be measured directly from remote-sensing observations, other remotely sensed variables can be used to estimate these fluxes. Moreover, these three fluxes are highly coupled and, therefore, a coupled model is ideal.

A fully connected neural network model was developed for this purpose [196], named water, energy, and carbon with artificial neural networks (WECANN). Inputs to WECANN are remotely sensed estimates of precipitation, soil moisture, net radiation, snow water equivalent, air temperature and solar induced fluorescence. The target variables for training the model were derived from outputs of global models. However, this presents the difficulty that the target variables are model outputs that can have substantial error, which will propagate in the WECANN model. To mitigate this problem, three independent estimates of each of the three fluxes (LE, H and GPP) were retrieved from the global models. Then a novel statistical approach, named triple collocation (TC, Figure S2, equation S1), was used to combine those estimates to a new dataset for training the WECANN model.

Triple collocation (TC) is a technique for estimating the unknown error (measured with standard deviations or RMSEs) of three mutually independent measurement systems, without treating any one system as zero-error “truth” [197]. The three measurement systems estimate a variable collocated in space and time, hence the name triple collocation. Using these probabilities, at each pixel and at each time one of the three estimates of the target variable is randomly selected to generate the TD.

The results of WECANN model outputs were evaluated against ground measurements from global FLUXNET towers from 2007 to 2015 (Figure 10), using both the coefficient of determination and RMSE to evaluate accuracy. These show that WECANN’s correlation was on average 17% higher (range 8–51%) than that of any one of the three individual inputs, while the RMSE was 21% lower (range 4–54%). These differences provide a partial quantification of the error inherent in any one of these training inputs and show that by combining them using the TC technique, we can reduce error in

an ML model for predicting the fluxes at global scale. This case study illustrates a means of assessing and accounting for error in TD for cases in which these data are not created specifically for the project, but rather are pre-existing data products with potentially quite different characteristics and potentially unknown error.

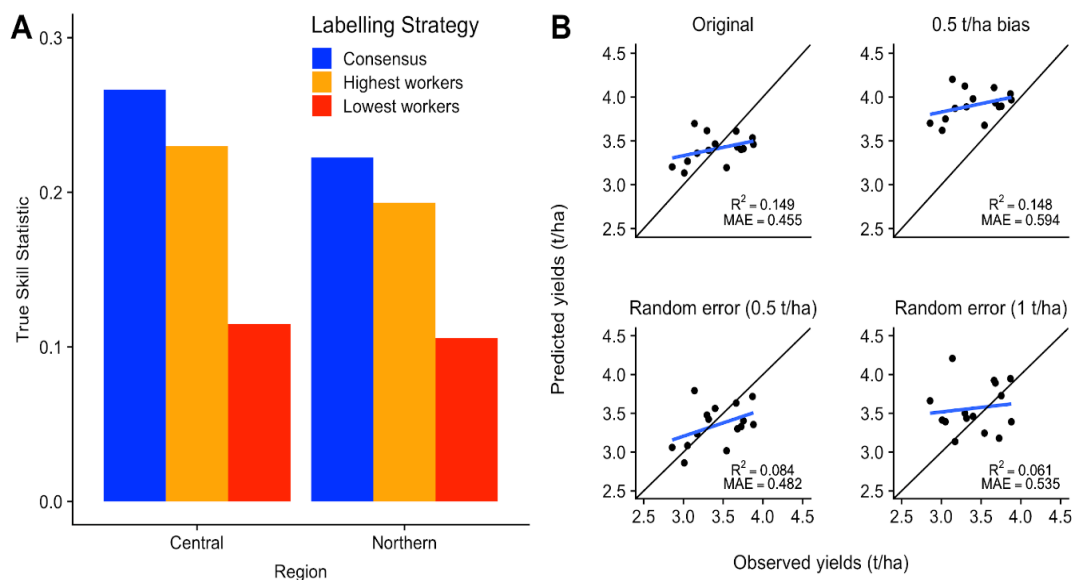


**Figure 10.** Coefficient of determination ( $R^2$ ) and root mean squared error (RMSE) of the water, energy, and carbon with artificial neural networks (WECANN) model output against ground measurements from FLUXNET towers in comparison to the three datasets used to generate the target training data for latent heat flux (LE) (a,b), sensible heat flux (H) (c,d) and gross primary production (GPP) (e,f).

### 3.3. Agricultural Monitoring

Two agricultural cases illustrate how TD error can impact both categorical and quantitative remotely sensed measures. The first relates to cropland mapping and is drawn from an ongoing study focused on mapping smallholder agricultural fields at high spatial resolution (3–4 m) in Ghana. The mapping method is based on “active learning”, in which a random forest-based [124,198,199] ML algorithm is iteratively trained and validated by a crowdsourcing platform. This platform enlists human trainers to visually interpret and digitize field boundaries in the imagery (PlanetScope visual and near-infrared surface reflectance [128]) being classified [149,150,198]. During this process, a protocol is used to assess the accuracy of training labels, in which each worker is periodically directed to a training reference site where the field boundaries are already known but are invisible to the worker. Using these training reference sites, the interpreters’ maps are then scored using a multi-dimensional accuracy

assessment algorithm [150], resulting in an average TD accuracy score for each worker ranging from 0 (complete disagreement with reference) to 1 (perfect agreement). Each label site is mapped by at least five workers, and the resulting worker-specific accuracy scores are used within a Bayesian merging algorithm to combine the five sets of labels into a single consensus label, which is then used to train the random forest classifier. Here we use the worker-specific training accuracy scores to assess the impact of label quality on map accuracy by assessing three variants of two random forest-generated maps, one over Central Ghana (~3400 km<sup>2</sup>) and one over Northern Ghana (~3100 km<sup>2</sup>). The first two maps were trained using labels generated by the worker with the least accurate TD, the second two by the most accurate worker, and the third using the consensus labels. The accuracy of each pair of maps was then assessed against the validation set (reserved consensus labels) using the true skill statistic [90] (sensitivity + specificity – 1, with scores ranging from –1 to 1). The results show a substantial difference in accuracy between the maps trained with the least and most accurate workers' labels (Figure 11A), with the former having 7–9% more skill than the latter, while maps based on consensus labels have ~3% more skill than those of the most accurate workers' labels.



**Figure 11.** A comparison (A) of the accuracy (based on the true skill statistic) of cropland maps over two areas of Ghana when generated by labels of different levels of quality (red = least accurate workers' labels; orange = most accurate workers' labels; blue = "consensus" labels made by merging all workers' labels). (B) Results from a random forest model of wheat yields trained on satellite-derived vegetation indices, showing the relationship between predicted yield and independent observed yields, in terms of the fit against the line and the regression slope of the relationship (points and regression line represent the mean of a single randomly selected model permutation). The average mean absolute error (MAE) and average regression  $R^2$ s calculated across all permutations are shown for each model.

The second case relates to remotely sensed crop estimates of wheat yields collected in 48 smallholder fields in Bihar, India in 2016–17 [200]. Yield data were collected via eight  $2 \times 1$  m<sup>2</sup> crop cuts within each field, and PlanetScope-derived green chlorophyll vegetation indices (GCVI) were calculated over each field from imagery collected over four dates during the growing season (13 January, 25 February, 12 March, and 14 April 2017). A random forest regression was trained on the yield measured for each field, using the four dates of GCVI values as predictors. To test the effect of TD error on the resulting yield predictions, three types of noise were artificially introduced into the yield data used for training: (1) a systematic 0.5 ton/ha overestimate with randomly distributed errors sampled from a normal distribution with a mean of 0 ton/ha, (2) random noise with standard deviations of 0.5 ton/ha, and (3) random noise with standard deviations of 1 ton/ha. A baseline model fit to unperturbed data was also developed. Each model was trained on three separate randomly selected subsets of

32 perturbed observations, and the predictions were made for the remaining 16 held-out (independent) yield observations, which were not perturbed. This three-fold cross validation process was repeated 50 times, with each permutation using a different random seed to construct the folds, in order to achieve stable error metrics. The model performance was assessed by calculating the averages of the mean absolute error (MAE) of the prediction, and the  $R^2$  of regressions fit between prediction and observed values (Figure 11B).

The results show that four models, including the baseline, compressed the range of yields, as seen in the shallow slope between observed versus predicted values, but prediction error was 18–31% higher when training yields had either the high level of random or systematic error within them. The smaller amount of random noise only added ~6% error to the predictions, suggesting that RandomForest is tolerant to some training error. Note that the average  $R^2$  of the observed-predicted regression fit was nearly the same for the systematic error case as the baseline, which shows that this metric can be an unreliable measure of performance for quantitative measures, and that it is important to assess fit against the  $y = x$  line and to use a metric such as mean absolute error.

#### 4. Guidelines and Recommendations

Our review and case studies show that the impacts of TD error on EO applications can vary, as can the procedures for assessing those impacts. Nevertheless, several best practices and guidelines can be discerned from this work. Below we synthesize a set of suggested steps for minimizing and accounting for TD error, within the context of undertaking and assessing the accuracy of a typical ML-based mapping project.

##### 4.1. Step 1: Define Acceptable Level of Accuracy and Choose Appropriate Metric

As a starting point, researchers should determine the minimum level of accuracy required for their application, using the accuracy metric(s) most appropriate for answering their questions [201]. For example, if the goal of creating a categorical map is to obtain an unbiased area estimate for a particular land cover, it is essential to account for the map's commission and omission errors by adjusting the proportional area estimate of the cover type derived from the map by the proportion of that type estimated from the map reference sample [39,40,57]. For a continuous variable in which the absolute accuracy of the mapped variable is most important, then the mean absolute deviation from the  $y = x$  line is more informative than  $R^2$  [93,94].

Error in the map reference data should also be factored into the selected accuracy metrics and resulting map-derived measures. Several published methods exist for categorical data (see Section 1.2.1). For continuous variables, the fit between the mapped and map reference variables can be assessed using Type 2 regression, which allows for error in the dependent (map reference) variable [202], unlike the more commonly used Type 1 regression. Determining map reference data error is critical to determining overall map accuracy. The error in these data effectively determines the upper limit of achievable map accuracy, as it is difficult (but not impossible; see [47]) to know whether a model's predictions are more accurate than its map reference data; thus if the map reference data are only 90% accurate, then the map can be at most 90% accurate. Acceptable accuracy should thus be determined relative to the accuracy of the map reference data, rather than the implicit assumption of 100%, which is widely used since map reference data are usually considered perfect [38,39,47,57,67].

Although the above steps relate primarily to concerns about map accuracy assessment, they are essential to establishing best practices for map TD. This is firstly due to the fact that, without undertaking rigorous accuracy assessment as described above, it is not possible to assess fully how TD error impacts map accuracy. And secondly, the processes of map reference data and TD generation are often tightly intertwined and impacted by many of the same sources of error (see Sections 1.2.1 and 1.2.2). The procedures for minimizing and measuring errors in both datasets are thus often the same. Our subsequent recommendations, therefore, cover both training and map reference datasets, except where we indicate necessary distinctions.

#### 4.2. Step 2: Minimize Design-Related Errors

The next logical step in a mapping project that relies on TD is to design strategies for independently collecting the training and map reference samples. Although there are numerous factors to consider, there are several general aspects of design that can help minimize potential TD errors.

##### 4.2.1. Sample Design

The first consideration relates to the sampling design itself, meaning where, when, how many, and what type of samples are placed (e.g., simple random, clustered, stratified, systematic, purposive/directed). With respect to the TD, this depends to a certain extent on the requirements of the selected ML algorithm, since various ML algorithms have differing requirements with respect to geographic distribution [53] and class balance of samples, e.g., [31,48,80]. Geographic representativeness and the degree to which the TD capture the variability in the feature of interest are also important TD sample design considerations [53,61,150,203]. Continuous TD, particularly those collected in situ, are often point samples. Therefore a sampling protocol should be used to match field measurements and pixel dimensions in order to avoid scaling problems associated with the modifiable areal unit problem [142,143].

The road mapping case study above shows the type of errors that can result when maps are trained with samples that do not adequately represent the features in a particular region. TD can in practice be highly localized or relevant for a limited spatial extent or temporal period [160,194]. This problem may continue to become more relevant, given the increase in stock or benchmark training libraries and subsequent attempts to transfer pre-trained models to different regions, time periods, or scales of observation [73,204]. While such benchmark libraries can be of immense benefit as TD for large area EO research, the representativeness of the features of interest should be assessed and augmented as needed, as shown above in the Khartoum model case study (Figure 9D). For some widely-used ML algorithms, such as random forests, the best practice appears to be to train with data collected within the mapping region (e.g., within a particular agroecoregion [55,205]), and to avoid over-generalizing or transferring models to other regions [206]. However, until more published studies are available, it is not clear whether this rule applies to deep-learning models. When using citizen science or crowdsourcing approaches to generate these data, representativeness can be ensured by directing labelers to the selected TD sites, e.g., [150], rather than having the interpreters select the regions to map.

Samples should also be temporally representative of the imagery that is being classified [61]. That is, relative to the imagery being classified, the TD (and map reference) sample should be collected within a window of time that matches the characteristic rate of change of the feature being mapped. This interval can be estimated by measuring the temporal autocorrelation in the feature of interest [207]. For rapidly changing phenomena, such as deforestation events, snow/ice melt, and vegetation coverage during phenological transition, the sample may need to be captured within a few days or weeks of the acquisition of the imagery being classified, whereas for slower-moving features a sample collected within a few years may be sufficient.

In cases where training and reference samples are generated simultaneously, it is essential that TD sample design does not undermine the standards required for an independent, probabilistic map reference sample, *sensu* [67]. This is particularly relevant for extremely large, geographically broad benchmark datasets, which may be used for either TD or map reference data, assuming the specific data set used conforms to the appropriate criteria. Stehman et al. [176] describe procedures for rigorous incorporation of crowdsourced data while maintaining an appropriate probability-based sampling approach specifically for map reference data, and Stehman and Foody [57] explore issues relating to independence between TD and map reference data during sample design. Beyond those considerations, it is important to note that the map reference sample's independence is compromised when it is used to iteratively refine the mapping algorithm. This problem can best be understood within the context of cross validation, which is appropriate for ML parameter tuning, e.g., [31]. However, when the number



of folds exceed one (as in our yield estimation case study; Figure 11B) then the portions excluded from training lose statistical independence and can no longer serve as the map reference [77]. Map reference data independence may also be undermined when training sites are selected iteratively, in order to increase their representativeness and improve ML performance e.g., [55,149]. If the gain due to new training sites is assessed against the map reference, then it will also lose independence after the first iteration. Moreover, any error in the map reference sample will be integrated into the final map. Xiong et al. [55] avoided this problem by visually assessing whether their classifier improved map quality after having new TD points added to the initial sample. A more quantitative approach is to divide an initial sample into three splits: one for training, the second for validating algorithm improvements, including those related to the addition of new training sites, and the third as the map reference, used only for final accuracy assessment. This partitioning approach can be implemented in the mapping platform used in the cropland mapping case study, Figure 11A [199].

#### 4.2.2. Training Data Sources

The requirements for temporal representativeness make the source of training imagery a critical consideration for projects that rely on image interpretation. The use of basemap imagery is not recommended for training maps of dynamic features, given their broad range and uneven distribution of image acquisition dates [61], unless the age of the imagery being classified can be matched to that of the training imagery. Otherwise, there is substantial potential for introducing error into the mapping algorithm (e.g., Figure 1), and its impact may be hard to assess, particularly if the map reference sample is collected from the basemap. The goal of temporal representativeness must be balanced with the need to have a sufficiently high spatial resolution to accurately interpret the smallest target features/classes (i.e., the MMU; see Step 3). Beyond matters of cost, this tradeoff is one reason that HR/VHR basemaps are widely used [61]. Newly available commercial imagery such as PlanetScope [128] are collected at high temporal frequency (near-daily) and have a spatial resolution sufficient for many visual interpretation tasks (3–4 m) and, therefore, may be a preferable source of training imagery for developing maps representing the post-2016 period. Finally, in designing an image-based sample, it is also important to consider additional characteristics that can influence interpreters' judgement, such as atmospheric quality (e.g., clouds, haze), sensor view angle, sun angle, spectral band selection, and image contrast stretches [74].

#### 4.2.3. Legend Design

For thematic maps, legend design merits special consideration as it relates to TD, particularly for multi-temporal and/or large area projects that rely on multiple image datasets [61]. As discussed in Section 2 above, objects of interest, including land-cover patches (i.e., the MMU), should be at least twice as large as the pixel resolution of the imagery used in the classification algorithm, assuming a requirement for spectrally pure pixels [136,168,208]. When image spatial resolution is too coarse relative to the scene elements of interest, image interpretation errors are likely due to mixed pixels [127,137,138]. This implies that in designing a legend, researchers should select classes with an MMU that is large enough to be effectively captured by the coarsest resolution imagery to be incorporated in the model, which will help avoid the problem of collecting training samples with mixed pixels [55]. This consideration is particularly relevant since HR/VHR imagery is often used to create TD and map reference data, while the mapping algorithm is applied to moderate- or coarse-resolution imagery, e.g., [55,120,209,210]. Alternatively, researchers may opt to select a classification workflow which explicitly incorporates mixed pixels, e.g., [98,165,173].

Spatial representativeness should be considered as a limiting factor for legend design [53], and to the extent possible, researchers should attempt to use categories that are supported by both the spatial resolution of the model data and the field sampling protocols to be used; we recommend that researchers consult the extensive literature on legend design [25,144–147,211–213].

#### 4.3. Step 3: Minimize Collection-Related Errors

There are numerous ways to collect TD for categorical and continuous mapping projects, each with their own sources of error. There are thus many potential approaches for minimizing the associated collection errors, which may be quite specific to a particular variable (e.g., for agricultural area estimates [214]). However, there are several general approaches that can be followed to minimize TD collection errors. Our focus here is primarily on error in image-interpreted TD, which is one of the most common approaches used to training ML mapping algorithms. We also touch on the specific case of model-derived training data.

Whenever possible, we recommend using protocols that incorporate training reference data to independently assess TD accuracy, particularly for image-interpreted TD, e.g., [150]. Training reference datasets can be limited in size compared to the ultimate sample size, provided that training reference locations are randomly presented to interpreters during the data creation campaign [150]. Active feedback during training label creation can also help reduce errors on a rolling basis by providing interpreters with information regarding their performance [174].

If comparison against training reference data is not possible, then consensus methods for generating TD may be the next best alternative. Consensus between at least 3 interpreters is recommended to allow for majority voting [34,46], but more complex land covers may require up to 7 interpreters [46]. Consensus among several domain experts may also be the best and most practical measure for collecting both training reference data and map reference data [34,57]. In the case of image-interpreted samples, consensus approaches should employ multiple interpreters to label the same site. For continuous variables, several independent or repeated in situ measurements should be made and aggregated. For modeled variables where the error is unknown, as in the surface flux case study, training based on the outputs of multiple independent models is recommended. The agricultural case study shows how multiple mappings can be used to quantify label uncertainty (Figure 12A) and minimize the amount of labeling error, yielding improved map accuracy (Figure 11A). The surface flux case study demonstrates these same benefits across several continuous variables (Figure 10). The number of separate measures or interpreters needed will vary depending on the application. For example, in the cropland mapping case study, 5 interpreters labeled each consensus training sample, and in the continuous surface flux example, 3 separate modeled inputs were used.

Further steps can be taken to minimize TD collection errors arising from image interpretation. Interpreters should be given thorough training regarding the task [34], which may include instruction on remote-sensing principles as well as local or regional contextual information. Local domain expertise is particularly helpful for consistent identification of idiosyncratic land covers [163]. Interpreter education is particularly important for crowdsourcing or citizen science data collection campaigns, as participants typically lack formal experience in image interpretation [151,215].

As described in Step 2 above, image interpretation is inadvisable when the available imagery does not support the legend categories in terms of spatial, spectral, temporal, or radiometric resolution [216–218]. Researchers must be especially cautious in the similar but potentially more hazardous case that HR/VHR imagery is used to create training samples that are then used with coarser resolution imagery when ingested into the ML model. Assuming that researchers correctly specify their data selection and legend design when using higher spatial resolution imagery to create TD, image interpretation errors due to insufficient resolution should be minimized; however, special care should be given to borderline classes, or classes exhibiting a high degree of spatial and/or spectral variability due to land-cover mixtures within the pixel [127,137,138,154,219]. In such cases, we recommend that training polygons be created near the center of scene objects, where pixel mixing is likely to be minimized, e.g., [55].

Another important error-minimizing approach relates to cases in which TD comes from a process model, as in the surface flux example outlined above. Process models are also increasingly used to train crop yield mapping models, due to the difficulty of obtaining sufficiently large and reliable field-scale yield data for training [220]. To circumvent this challenge, the scalable yield mapping (SCYM) method [221,222] uses a mechanistic crop model to simulate yields under various environmental and management conditions. The model's outputs then become inputs for training an empirical mapping model (typically ML), in which the simulated yield is the dependent variable and a subset of remotely retrievable model variables serve as predictors. TD errors in such cases can be minimized by rigorously calibrating the process model (itself a challenging task) using best practices from the relevant modeling literature, e.g., [223]. Alternatively, if modeled TD are necessary but careful calibration is not possible (e.g., because the data are pre-existing), then a merging approach such as triple collocation (Section 3.2) can help reduce training error.

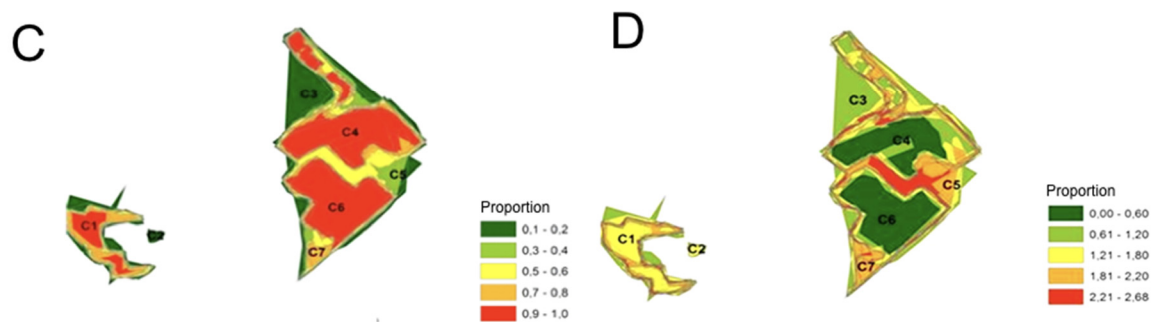
#### 4.4. Step 4. Assess Error in Training Data Error

The best way to assess both TD (and map reference data) error is to measure it directly. For continuous variables, calculating measurement error should be possible in many cases, even for model-generated data, in which the variance can be calculated from simulation treatments, e.g., [223]. For categorical mapping, label error can be measured using an internal accuracy assessment protocol that makes use of predefined training reference data (e.g., Estes et al., [150]).

However, it can be challenging to produce training reference data, and indeed in some cases the true category is not clear, whether looking at an image or standing on site. In these cases, or when a direct TD error measurement protocol is not available, we recommend that researchers calculate uncertainty estimates based on repeated measures or multiple interpreter approaches (e.g., the crowd standard deviation [151]) described in Step 3 above (and see Figure 12); this is useful for both training and map reference data. We also recommend that additional measures relating to data collection speed, precision, and consistency be collected for individual data creators, as these can generate further insight into relative TD errors. This recommendation is based on experience in crowdsourced data creation [150,151], but it is applicable to any type of data collection, and could greatly bolster the understanding and quantification of error propagation. If it is not possible to either directly quantify TD error or relative uncertainty, then researchers should at a minimum clearly document the data creation methods, and detail likely sources of error and potential uncertainties.



Figure 12. Cont.



**Figure 12.** Two examples of consensus-based mapping approaches and their potential use for assessing training (or reference) data uncertainty. Panel (A) shows a collection of crop field boundary polygons drawn by five independent workers around crop fields visible in PlanetScope imagery collected over Ghana. These labels can be converted into a heat map (B) showing the overall agreement, the inverse of uncertainty. Similarly, 19 independent experts were asked to delineate slum settlements in image subset from Cape Town, South Africa. The polygons are converted into overall agreement and the uncertainty is modeled using random sets (C) shows the covering function, which is then used to calculate standard deviation of random set (D). Both these metrics indicate the variability as well as stability in boundaries delineated by different experts. Adapted with permission from Kohli et al. [163].

#### 4.5. Step 5. Evaluate and Communicate the Impact of Training Data Error

##### 4.5.1. TD Treatment Tiers

Due to the wide range of remote-sensing research currently underway, a wide variety of TD and classification algorithms are in use. Therefore, it is not possible to specify a single protocol for treatment of TD error. Instead, we outline three tiers that represent different levels of accounting for the impact of TD errors on resulting map products. These three tiers presuppose that researchers follow best practices for map accuracy assessment, which includes selecting the most appropriate, literature-recommended accuracy measure(s), quantifying map reference sample error, and accounting for the impact of map reference data error on the accuracy measures (per Step 1). If these best practices are followed, TD error impacts will already be implicitly accounted for within the accuracy measures, and the selected TD accounting tier will be governed by the purposes of the mapping application.

##### Tier 1

The optimal TD accuracy assessment, termed Tier 1, involves quantifying TD error using gold standard training reference data (Step 4). This information is then used to quantify various characteristics of the TD sample such as class balance and sample size. It is also used to determine the impacts of collection error stemming from label or measurement errors on model uncertainty and map accuracy (see Sections 1.2.2 and 2.2). For example, the impact of TD error on the certainty of random forest classifications can be assessed using measures derived from the margin function [48]. The impact of TD error on map accuracy should also be assessed by training models with TD adjusted to reflect the range of measured TD error, as illustrated by our cropland mapping case study, and with respect to variations in TD sample size and class balance [30,48,149]. This approach can be used to inform the researcher how much map improvement can be obtained by improving TD quality. As such, these tests should be performed against the validation sample rather than the map reference data, in order to preserve the independence of the latter.

We recommend that developers of benchmark TD libraries adhere to the tier 1 guidelines, keeping in mind that these data sets are likely to be used for a variety of purposes, including as TD and map reference data. Undertaking such evaluations can provide users important information about appropriate usage of these data for different ML models and geographies, and whether the benchmark data are appropriate for use as TD, training reference data, validation data, and/or map reference data. A rigorous quantification of error in the samples themselves is particularly important, since such

data are often likely to be used as training and/or map reference data. We strongly urge researchers to consider what purposes these benchmark data sets are appropriate for, and refer the reader to previously published literature regarding incorporation of non-probabilistic samples [176]. Ideally, this tier should also be followed by the makers of map products intended for widespread public use, who should also release TD and map reference data that were used during map creation [57]. This step would allow users full insight into the quality and usability of the map for their own purposes.

Published TD (and map reference data) should be documented with standard metadata, as shown in Table S4, including the relevant error metric associated with each observation. The SpatioTemporal Asset Catalog (STAC, <https://stacspec.org/>) provides a framework for standardization of metadata for EO data and is increasingly seen as an international standard for geospatial data.

#### Tier 2

If it is not possible to directly measure and quantify TD error, the next best approach to account for TD error is to introduce a plausible range of simulated error into the TD and evaluate its impact on model uncertainty and map accuracy after training separate models with the perturbed datasets [48]. If multiple workers are tasked with collecting TD for the same site, then the variance in their data can be calculated [151] to derive the uncertainty bounds (e.g., Figure 12). This approach is demonstrated in the building mapping case study (Section 3.1.1), which illustrates the sensitivity of key accuracy metrics to two different kinds of simulated labeling errors. The wheat yield case study (see Section 3.3) provides an example of this approach for a continuous variable.

This tier may also provide an acceptable standard for both benchmark datasets and publicly released map products, particularly where absolute error quantification is less important, as well as for publicly released map products. TD and map reference data should also be made openly available with standard metadata, as described above, including the uncertainty metric for each observation. If it is not possible to publish them (e.g., because of privacy concerns), then researchers should provide documentation that summarizes these data and their uncertainty.

#### Tier 3

If the TD error quantification in Tiers 1 or 2 are not possible, then researchers should at minimum publish their TD and map reference data, e.g., [55] with accompanying metadata that includes descriptions of potential errors and uncertainties. If data cannot be made openly available, then researchers should publish full descriptions of the potential error in the data. Adherence to this tier, at least the reporting component, should be the minimal standard practice in peer-reviewed, map-based scientific research.

#### 4.5.2. Communicating Error

Finally, uncertainty in ML-generated maps associated with both TD and map reference error should be faithfully reported within the maps and accompanying documents. Incomplete error reporting serves to limit the scientific validity and usefulness of these products [57]. Given that ML-generated maps are increasingly used by the public and policy domains, we advise makers of widely used maps to communicate these uncertainties and their consequences in a manner that is clear and understandable for broad audiences, including non-specialists, so that users can understand the map and its limitations. In general, we recommend including the error on or very close to the actual map, whether by means of metrics, the error matrix, and/or by using cartographic techniques for representing uncertainty. Examples of effective cartographic techniques for conveying uncertainty include selection of appropriate, intuitive, and color-blind friendly color schemes for classes and symbols, varying color value and saturation and font/line weight to indicate levels of uncertainty, use of crisp versus blurred boundaries and symbols to indicate the range of uncertainty, or display of consensus maps or side-by-side juxtaposition in cases of multiple, mutually exclusive predictions for

the same place and time (e.g., representing differently specified models) [42,43]. Maps of consensus in training labels can provide valuable uncertainty information to users, such as shown in Figure 12A,B.

#### 4.5.3. Towards an Open Training Data Repository

For the scientific community, the ideal standard of openness and replicability is to provide a complete description of TD collection practices, appropriate accuracy metrics, and perhaps most importantly of all, the raw data. We recommend the creation of a centralized, open source database of all available and relevant TD, using the details collected in the proposed template (Table S4), and recorded using the STAC framework. This type of open repository, taking inspiration from similar large-scale databases for computer vision (ImageNet, SIFT10M Dataset [224,225], and remote sensing (BigEarthNet, DeepSat, UC Merced Land-Use Dataset [73,226,227], should contain full training metadata, citations to the peer-reviewed literature, as well as links to downloadable versions of TD collection protocols. Following the philosophy of free and open source software, we strongly recommend that researchers embrace open source data, which is the only way by which a study can be truly reproduced.

## 5. Conclusions

Current practices in EO research are generally inattentive to the need to evaluate and communicate the impact of TD error on ML-generated maps. This oversight undermines the goals of scientific reproducibility and may compromise the insights drawn from the resulting maps. Improving these practices is important due to the increasing use of TD-intensive ML algorithms, which have motivated our review and recommendations.

To resolve terminological differences arising from the influence of non-EO disciplines, and to help contextualize TD considerations relative to established map accuracy assessment practice, we distinguish between four types of “truth” data used in ML-based mapping projects (training, validation, training reference, and map reference data), and define the appropriate role for each (Section 1.2). We identify causes of error in TD as well as map reference data, distinguishing where these vary (Section 2.1). We then explore the impacts of TD error (Section 2.2) and provide a set of case studies to illustrate the consequences of such error across a range of ML-based mapping applications (Section 3).

We then provide a set of guidelines for minimizing error arising from the design and collection of TD samples, and present recommendations for measuring and accounting for the impact of these errors (Section 4). Many of these guidelines and procedures also relate to map reference data generation, and we ground our recommendations in the existing best practices for map accuracy assignment (Sections 1.2.1 and 4.1). We conclude by defining three tiers of TD error accounting and reporting standards, which are designed to accommodate a wide range of ML-based mapping projects. The highest tiers should be adopted when creating open training libraries and public map products, both of which are increasingly being developed to meet the growing demand for EO-derived maps. In this context, there is a pressing need to rigorously evaluate the training requirements and relative performance of deep-learning models as they become more widely used for EO [36]. While TD is more visible in the context of LCLU and other categorical mapping projects, the need for rigorous, well-documented TD is also critically important for continuous variable applications in Earth System Sciences (e.g., hydrological research [228]). If adopted within the peer-reviewed literature, the standards we propose for TD treatment may improve confidence in scientific findings drawn from map-based research, which can otherwise be confounded by poorly quantified map errors [33,57].

**Supplementary Materials:** The following are available online at <http://www.mdpi.com/2072-4292/12/6/1034/s1>, Figure S1: Sample prediction results in Kumasi, Ghana. (a) Input imagery. (b) Predictions from the Las Vegas model. (c) Predictions from the Khartoum model. (d) Prediction from the Kumasi model. (e) Predictions from the Khartoum Model retrained in Kumasi. Figure S2: Schematic of product selection using the Triple Collocation approach. Table S1: List of peer-reviewed publications retrieved using Google Scholar search algorithm results. Table S2: Summary of commonly used error metrics. Table S3: Quantitative results of comparing each of the three models trained for the road detection case in Kumasi, Ghana to the validation labels. Table S4: Template and procedure for documenting training data.

**Author Contributions:** This article synthesizes the ideas of the 20 authors resulting from a workshop focused on issues of error in training data for Machine Learning approaches in Earth Observation research. Conceptualization, L.E. and A.E.; formal analysis and investigation, L.E., H.A., R.A., J.R.E., L.F., D.K., D.L., A.B.R., L.S., S.Y., Z.-F.Y.; writing—original draft preparation, A.E., L.E., H.A., L.F., D.K., D.L., R.G.P.J., A.B.R., Z.-F.Y.; writing—review and editing, A.E., L.E., K.C., H.A., R.A., L.F., M.A.F., M.J., D.K., J.C.L.B., J.L.M., J.R.; visualization, H.A., R.A., L.F., D.K., A.B.R., H.S.; supervision, L.E.; project administration, L.E. and A.E.; funding acquisition, L.E., A.B.R. All authors have read and agreed to the published version of the manuscript.

**Acknowledgments:** This work represents a synthesis of findings from a workshop held at Clark University on 8–9 January 2019. The workshop and subsequent paper writing and development was supported by a grant from Omidyar Network’s Property Rights Initiative, now PlaceFund. Additional support for developing methods and data presented here was provided by NASA (80NSSC18K0158), the National Science Foundation (SES-1801251), National Institute of Standards and Technology (2017-67003-26615), National Institute of Standards and Technology Summer Undergraduate Research Fellowship Program, and New York State Department of Environmental Conservation (DEC01-T00640GG-3350000). We thank Victoria Gammino for helpful input and advice, and David Allen, Ayo Deas, Lucy Hutyra, Clare Kohler, Barry Logan, Jaret Reblin, Ian Smith for assistance with fieldwork and data compilation.

**Conflicts of Interest:** The authors declare no conflict of interest.

## References

- Chen, J.; Chen, J.; Liao, A.; Cao, X.; Chen, L.; Chen, X.; He, C.; Han, G.; Peng, S.; Lu, M.; et al. Global Land Cover Mapping at 30 m Resolution: A POK-Based Operational Approach. *ISPRS J. Photogramm. Remote Sens.* **2015**, *103*, 7–27. [[CrossRef](#)]
- Friedl, M.A.; Sulla-Menashe, D.; Tan, B.; Schneider, A.; Ramankutty, N.; Sibley, A.; Huang, X. MODIS Collection 5 global land cover: Algorithm refinements and characterization of new datasets. *Remote Sens. Environ.* **2010**, *114*, 168–182. [[CrossRef](#)]
- Song, X.-P.; Hansen, M.C.; Stehman, S.V.; Potapov, P.V.; Tyukavina, A.; Vermote, E.F.; Townshend, J.R. Global land change from 1982 to 2016. *Nature* **2018**, *560*, 639–643. [[CrossRef](#)] [[PubMed](#)]
- Mohanty, B.P.; Cosh, M.H.; Lakshmi, V.; Montzka, C. Soil Moisture Remote Sensing: State-of-the-Science. *Vadose Zone J.* **2017**, *16*. [[CrossRef](#)]
- Daudt, R.C.; Le Saux, B.; Boulch, A.; Gousseau, Y. Guided Anisotropic Diffusion and Iterative Learning for Weakly Supervised Change Detection. *arXiv* **2019**.
- Hecht, R.; Meinel, G.; Buchroithner, M. Automatic identification of building types based on topographic databases—A comparison of different data sources. *Int. J. Cartogr.* **2015**, *1*, 18–31. [[CrossRef](#)]
- Zhang, X.; Jayavelu, S.; Liu, L.; Friedl, M.A.; Henebry, G.M.; Liu, Y.; Schaaf, C.B.; Richardson, A.D.; Gray, J. Evaluation of land surface phenology from VIIRS data using time series of PhenoCam imagery. *Agric. For. Meteorol.* **2018**, *256–257*, 137–149. [[CrossRef](#)]
- Tan, B.; Morissette, J.T.; Wolfe, R.E.; Gao, F.; Ederer, G.A.; Nightingale, J.; Pedelty, J.A. An Enhanced TIMESAT Algorithm for Estimating Vegetation Phenology Metrics From MODIS Data. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.* **2011**, *4*, 361–371. [[CrossRef](#)]
- Zhang, X.; Friedl, M.A.; Schaaf, C.B. Global vegetation phenology from Moderate Resolution Imaging Spectroradiometer (MODIS): Evaluation of global patterns and comparison with in situ measurements: GLOBAL PHENOLOGY FROM MODIS. *J. Geophys. Res.* **2006**, *111*, 981. [[CrossRef](#)]
- Schaaf, C.B.; Gao, F.; Strahler, A.H.; Lucht, W.; Li, X.; Tsang, T.; Strugnell, N.C.; Zhang, X.; Jin, Y.; Muller, J.-P.; et al. First operational BRDF, albedo nadir reflectance products from MODIS. *Remote Sens. Environ.* **2002**, *83*, 135–148. [[CrossRef](#)]
- Liu, Y.; Wang, Z.; Sun, Q.; Erb, A.M.; Li, Z.; Schaaf, C.B.; Zhang, X.; Román, M.O.; Scott, R.L.; Zhang, Q.; et al. Evaluation of the VIIRS BRDF, Albedo and NBAR products suite and an assessment of continuity with the long term MODIS record. *Remote Sens. Environ.* **2017**, *201*, 256–274. [[CrossRef](#)]
- Wang, Z.; Schaaf, C.B.; Sun, Q.; Shuai, Y.; Román, M.O. Capturing rapid land surface dynamics with Collection V006 MODIS BRDF/NBAR/Albedo (MCD43) products. *Remote Sens. Environ.* **2018**, *207*, 50–64. [[CrossRef](#)]
- Wan, Z. New refinements and validation of the MODIS Land-Surface Temperature/Emissivity products. *Remote Sens. Environ.* **2008**, *112*, 59–74. [[CrossRef](#)]

14. Jiménez-Muñoz, J.C.; Sobrino, J.A.; Skoković, D.; Mattar, C.; Cristóbal, J. Land Surface Temperature Retrieval Methods From Landsat-8 Thermal Infrared Sensor Data. *IEEE Geosci. Remote Sens. Lett.* **2014**, *11*, 1840–1843. [[CrossRef](#)]
15. Jean, N.; Burke, M.; Xie, M.; Davis, W.M.; Lobell, D.B.; Ermon, S. Combining satellite imagery and machine learning to predict poverty. *Science* **2016**, *353*, 790–794. [[CrossRef](#)] [[PubMed](#)]
16. Pekel, J.-F.; Cottam, A.; Gorelick, N.; Belward, A.S. High-resolution mapping of global surface water and its long-term changes. *Nature* **2016**, *540*, 418–422. [[CrossRef](#)]
17. Hansen, M.C.; Potapov, P.; Tyukavina, A. Comment on “Tropical forests are a net carbon source based on aboveground measurements of gain and loss”. *Science* **2019**, *363*. [[CrossRef](#)]
18. Gutierrez-Velez, V.H.; Pontius, R.G. Influence of carbon mapping and land change modelling on the prediction of carbon emissions from deforestation. *Environ. Conserv.* **2012**, *39*, 325–336. [[CrossRef](#)]
19. Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In Proceedings of the 2009 IEEE Conference on Computer Vision and Pattern Recognition, Miami, FL, USA, 20–25 June 2009; pp. 248–255.
20. Helber, P.; Bischke, B.; Dengel, A.; Borth, D. EuroSAT: A Novel Dataset and Deep Learning Benchmark for Land Use and Land Cover Classification. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens* **2019**, 1–10. [[CrossRef](#)]
21. Liu, Q.; Hang, R.; Song, H.; Li, Z. Learning Multiscale Deep Features for High-Resolution Satellite Image Scene Classification. *IEEE Trans. Geosci. Remote Sens.* **2018**, *56*, 117–126. [[CrossRef](#)]
22. Laso Bayas, J.C.; Lesiv, M.; Waldner, F.; Schucknecht, A.; Duerauer, M.; See, L.; Fritz, S.; Fraisl, D.; Moorthy, I.; McCallum, I.; et al. A global reference database of crowdsourced cropland data collected using the Geo-Wiki platform. *Sci. Data* **2017**, *4*, 170136. [[CrossRef](#)] [[PubMed](#)]
23. Lary, D.J.; Zewdie, G.K.; Liu, X.; Wu, D.; Levetin, E.; Allee, R.J.; Malakar, N.; Walker, A.; Mussa, H.; Mannino, A.; et al. Machine Learning Applications for Earth Observation. In *Earth Observation Open Science and Innovation*; Mathieu, P.-P., Aubrecht, C., Eds.; Springer International Publishing: Cham, Switzerland, 2018; pp. 165–218. ISBN 9783319656335.
24. Lary, D.J.; Alavi, A.H.; Gandomi, A.H.; Walker, A.L. Machine learning in geosciences and remote sensing. *Geosci. Front.* **2016**, *7*, 3–10. [[CrossRef](#)]
25. Loveland, T.R.; Reed, B.C.; Brown, J.F.; Ohlen, D.O.; Zhu, Z.; Yang, L.; Merchant, J.W. Development of a global land cover characteristics database and IGBP DISCover from 1 km AVHRR data. *Int. J. Remote Sens.* **2000**, *21*, 1303–1330. [[CrossRef](#)]
26. Sulla-Menashe, D.; Gray, J.M.; Abercrombie, S.P.; Friedl, M.A. Hierarchical mapping of annual global land cover 2001 to present: The MODIS Collection 6 Land Cover product. *Remote Sens. Environ.* **2019**, *222*, 183–194. [[CrossRef](#)]
27. Fortier, J.; Rogan, J.; Woodcock, C.E.; Runfola, D.M. Utilizing Temporally Invariant Calibration Sites to Classify Multiple Dates and Types of Satellite Imagery. *Photogramm. Eng. Remote Sens.* **2011**, *77*, 181–189. [[CrossRef](#)]
28. Foody, G.M.; Mathur, A. Toward intelligent training of supervised image classifications: Directing training data acquisition for SVM classification. *Remote Sens. Environ.* **2004**, *93*, 107–117. [[CrossRef](#)]
29. Graves, S.J.; Asner, G.P.; Martin, R.E.; Anderson, C.B.; Colgan, M.S.; Kalantari, L.; Bohlman, S.A. Tree Species Abundance Predictions in a Tropical Agricultural Landscape with a Supervised Classification Model and Imbalanced Data. *Remote Sens.* **2016**, *8*, 161. [[CrossRef](#)]
30. Foody, G.; Pal, M.; Rocchini, D.; Garzon-Lopez, C. The sensitivity of mapping methods to reference data quality: Training supervised image classifications with imperfect reference data. *Int. J. Geo-Inf.* **2016**, *5*, 199. [[CrossRef](#)]
31. Maxwell, A.E.; Warner, T.A.; Fang, F. Implementation of machine-learning classification in remote sensing: an applied review. *Int. J. Remote Sens.* **2018**, *39*, 2784–2817. [[CrossRef](#)]
32. Huang, C.; Davis, L.S.; Townshend, J.R.G. An assessment of support vector machines for land cover classification. *Int. J. Remote Sens.* **2002**, *23*, 725–749. [[CrossRef](#)]
33. Estes, L.; Chen, P.; Debats, S.; Evans, T.; Ferreira, S.; Kuemmerle, T.; Ragazzo, G.; Sheffield, J.; Wolf, A.; Wood, E.; et al. A large-area, spatially continuous assessment of land cover map error and its impact on downstream analyses. *Glob. Chang. Biol.* **2018**, *24*, 322–337. [[CrossRef](#)] [[PubMed](#)]



34. Pengra, B.W.; Stehman, S.V.; Horton, J.A.; Dockter, D.J.; Schroeder, T.A.; Yang, Z.; Cohen, W.B.; Healey, S.P.; Loveland, T.R. Quality control and assessment of interpreter consistency of annual land cover reference data in an operational national monitoring program. *Remote Sens. Environ.* **2019**, *111*, 261. [[CrossRef](#)]
35. Zhu, X.X.; Tuia, D.; Mou, L.; Xia, G.; Zhang, L.; Xu, F.; Fraundorfer, F. Deep Learning in Remote Sensing: A Comprehensive Review and List of Resources. *IEEE Geosci. Remote Sens. Mag.* **2017**, *5*, 8–36. [[CrossRef](#)]
36. Ma, L.; Liu, Y.; Zhang, X.; Ye, Y.; Yin, G.; Johnson, B.A. Deep learning in remote sensing applications: A meta-analysis and review. *ISPRS J. Photogramm. Remote Sens.* **2019**, *152*, 166–177. [[CrossRef](#)]
37. Foody, G.M. Status of land cover classification accuracy assessment. *Remote Sens. Environ.* **2002**, *80*, 185–201. [[CrossRef](#)]
38. Foody, G.M. Assessing the accuracy of land cover change with imperfect ground reference data. *Remote Sens. Environ.* **2010**, *114*, 2271–2285. [[CrossRef](#)]
39. Olofsson, P.; Foody, G.M.; Herold, M.; Stehman, S.V.; Woodcock, C.E.; Wulder, M.A. Good practices for estimating area and assessing accuracy of land change. *Remote Sens. Environ.* **2014**, *148*, 42–57. [[CrossRef](#)]
40. Pontius, R.G.; Millones, M. Death to Kappa: Birth of quantity disagreement and allocation disagreement for accuracy assessment. *Int. J. Remote Sens.* **2011**, *32*, 4407–4429. [[CrossRef](#)]
41. Congalton, R.G.; Green, K. *Assessing the Accuracy of Remotely Sensed Data: Principles and Practices*; CRC Press: Boca Raton, FL, USA, 2008.
42. Monmonier, M. Cartography: Uncertainty, interventions, and dynamic display. *Prog. Hum. Geogr.* **2006**, *30*, 373–381. [[CrossRef](#)]
43. MacEachren, A.M. Visualizing Uncertain Information. *Cartogr. Perspect.* **1992**, *1*, 10–19. [[CrossRef](#)]
44. Goodchild, M.F.; Gopal, S. *The Accuracy of Spatial Databases*; CRC Press: Boca Raton, FL, USA, 1989; ISBN 9780203490235.
45. Congalton, R.G. A review of assessing the accuracy of classifications of remotely sensed data. *Remote Sens. Environ.* **1991**, *37*, 35–46. [[CrossRef](#)]
46. McRoberts, R.E.; Stehman, S.V.; Liknes, G.C.; Næsset, E.; Sannier, C.; Walters, B.F. The effects of imperfect reference data on remote sensing-assisted estimators of land cover class proportions. *ISPRS J. Photogramm. Remote Sens.* **2018**, *142*, 292–300. [[CrossRef](#)]
47. Carlotto, M.J. Effect of errors in ground truth on classification accuracy. *Int. J. Remote Sens.* **2009**, *30*, 4831–4849. [[CrossRef](#)]
48. Mellor, A.; Boukir, S.; Haywood, A.; Jones, S. Exploring issues of training data imbalance and mislabelling on random forest performance for large area land cover classification using the ensemble margin. *ISPRS J. Photogramm. Remote Sens.* **2015**, *105*, 155–168. [[CrossRef](#)]
49. Swan, B.; Laverdiere, M.; Yang, H.L. How Good is Good Enough?: Quantifying the Effects of Training Set Quality. In Proceedings of the 2Nd ACM SIGSPATIAL International Workshop on AI for Geographic Knowledge Discovery, Seattle, WA, USA, 6 November 2018; ACM: New York, NY, USA, 2018; pp. 47–51.
50. Ghimire, B.; Rogan, J.; Galiano, V.R.; Panday, P.; Neeti, N. An Evaluation of Bagging, Boosting, and Random Forests for Land-Cover Classification in Cape Cod, Massachusetts, USA. *GISci. Remote Sens.* **2012**, *49*, 623–643. [[CrossRef](#)]
51. Rodriguez-Galiano, V.F.; Ghimire, B.; Rogan, J.; Chica-Olmo, M.; Rigol-Sanchez, J.P. An assessment of the effectiveness of a random forest classifier for land-cover classification. *ISPRS J. Photogramm. Remote Sens.* **2012**, *67*, 93–104. [[CrossRef](#)]
52. Bruzzone, L.; Persello, C. A Novel Context-Sensitive Semisupervised SVM Classifier Robust to Mislabeled Training Samples. *IEEE Trans. Geosci. Remote Sens.* **2009**, *47*, 2142–2154. [[CrossRef](#)]
53. Cracknell, M.J.; Reading, A.M. Geological mapping using remote sensing data: A comparison of five machine learning algorithms, their response to variations in the spatial distribution of training data and the use of explicit spatial information. *Comput. Geosci.* **2014**, *63*, 22–33. [[CrossRef](#)]
54. Mellor, A.; Boukir, S. Exploring diversity in ensemble classification: Applications in large area land cover mapping. *ISPRS J. Photogramm. Remote Sens.* **2017**, *129*, 151–161. [[CrossRef](#)]
55. Xiong, J.; Thenkabail, P.S.; Tilton, J.C.; Gumma, M.K.; Teluguntla, P.; Oliphant, A.; Congalton, R.G.; Yadav, K.; Gorelick, N. Nominal 30-m Cropland Extent Map of Continental Africa by Integrating Pixel-Based and Object-Based Algorithms Using Sentinel-2 and Landsat-8 Data on Google Earth Engine. *Remote Sens.* **2017**, *9*, 1065. [[CrossRef](#)]

56. Bey, A.; Jetimane, J.; Lisboa, S.N.; Ribeiro, N.; Siteo, A.; Meyfroidt, P. Mapping smallholder and large-scale cropland dynamics with a flexible classification system and pixel-based composites in an emerging frontier of Mozambique. *Remote Sens. Environ.* **2020**, *239*, 111611. [[CrossRef](#)]
57. Stehman, S.V.; Foody, G.M. Key issues in rigorous accuracy assessment of land cover products. *Remote Sens. Environ.* **2019**, *231*, 111199. [[CrossRef](#)]
58. Zhang, C.; Xie, Z. Object-based Vegetation Mapping in the Kissimmee River Watershed Using HyMap Data and Machine Learning Techniques. *Wetlands* **2013**, *33*, 233–244. [[CrossRef](#)]
59. Rogan, J.; Franklin, J.; Stow, D.; Miller, J.; Woodcock, C.; Roberts, D. Mapping land-cover modifications over large areas: A comparison of machine learning algorithms. *Remote Sens. Environ.* **2008**, *112*, 2272–2283. [[CrossRef](#)]
60. Copass, C.; Antonova, N.; Kennedy, R. Comparison of Office and Field Techniques for Validating Landscape Change Classification in Pacific Northwest National Parks. *Remote Sens.* **2018**, *11*, 3. [[CrossRef](#)]
61. Lesiv, M.; See, L.; Laso Bayas, J.C.; Sturn, T.; Schepaschenko, D.; Karner, M.; Moorthy, I.; McCallum, I.; Fritz, S. Characterizing the Spatial and Temporal Availability of Very High Resolution Satellite Imagery in Google Earth and Microsoft Bing Maps as a Source of Reference Data. *Land* **2018**, *7*, 118. [[CrossRef](#)]
62. Biradar, C.M.; Thenkabail, P.S.; Noojipady, P.; Li, Y.; Dheeravath, V.; Turrall, H.; Velpuri, M.; Gumma, M.K.; Gangalakunta, O.R.P.; Cai, X.L.; et al. A global map of rainfed cropland areas (GMRCA) at the end of last millennium using remote sensing. *Int. J. Appl. Earth Obs. Geoinf.* **2009**, *11*, 114–129. [[CrossRef](#)]
63. Mallinis, G.; Emmanoloudis, D.; Giannakopoulos, V.; Maris, F.; Koutsias, N. Mapping and interpreting historical land cover/land use changes in a Natura 2000 site using earth observational data: The case of Nestos delta, Greece. *Appl. Geogr.* **2011**, *31*, 312–320. [[CrossRef](#)]
64. Jawak, S.D.; Luis, A.J. Improved land cover mapping using high resolution multiangle 8-band WorldView-2 satellite remote sensing data. *JARS* **2013**, *7*, 073573. [[CrossRef](#)]
65. Ye, S.; Pontius, R.G.; Rakshit, R. A review of accuracy assessment for object-based image analysis: From per-pixel to per-polygon approaches. *ISPRS J. Photogramm. Remote Sens.* **2018**, *141*, 137–147. [[CrossRef](#)]
66. Fritz, S.; See, L.; Perger, C.; McCallum, I.; Schill, C.; Schepaschenko, D.; Duerauer, M.; Karner, M.; Dresel, C.; Laso-Bayas, J.-C.; et al. A global dataset of crowdsourced land cover and land use reference data. *Sci. Data* **2017**, *4*, 170075. [[CrossRef](#)] [[PubMed](#)]
67. Stehman, S.V. Sampling designs for accuracy assessment of land cover. *Int. J. Remote Sens.* **2009**, *30*, 5243–5272. [[CrossRef](#)]
68. Brodrick, P.G.; Davies, A.B.; Asner, G.P. Uncovering Ecological Patterns with Convolutional Neural Networks. *Trends Ecol. Evol.* **2019**, *34*, 734–745. [[CrossRef](#)] [[PubMed](#)]
69. Xiao, T.; Xia, T.; Yang, Y.; Huang, C.; Wang, X. Learning from massive noisy labeled data for image classification. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Boston, MA, USA, 7–12 June 2015; pp. 2691–2699.
70. Frénay, B.; Verleysen, M. Classification in the presence of label noise: A survey. *IEEE Trans. Neural Netw. Learn. Syst.* **2014**, *25*, 845–869.
71. Brodley, C.E.; Friedl, M.A. Identifying Mislabeled Training Data. *J. Artif. Intell. Res.* **1999**, *11*, 131–167. [[CrossRef](#)]
72. Van Etten, A.; Lindenbaum, D.; Bacastow, T.M. SpaceNet: A Remote Sensing Dataset and Challenge Series. *arXiv* **2018**.
73. Sumbul, G.; Charfuelan, M.; Demir, B.; Markl, V. BigEarthNet: A Large-Scale Benchmark Archive For Remote Sensing Image Understanding. *arXiv* **2019**.
74. Lesiv, M.; Laso Bayas, J.C.; See, L.; Duerauer, M.; Dahlia, D.; Durando, N.; Hazarika, R.; Kumar Sahariah, P.; Vakolyuk, M.; Blyshchyk, V.; et al. Estimating the global distribution of field size using crowdsourcing. *Glob. Chang. Biol.* **2019**, *25*, 174–186. [[CrossRef](#)]
75. Fritz, S.; McCallum, I.; Schill, C.; Perger, C.; See, L.; Schepaschenko, D.; van der Velde, M.; Kraxner, F.; Obersteiner, M. Geo-Wiki: An Online Platform for Improving Global Land Cover. *Environ. Model. Softw.* **2012**, *31*, 110–123.
76. Goodchild, M.F. Citizens as sensors: The world of volunteered geography. *GeoJournal* **2007**, *69*, 211–221. [[CrossRef](#)]
77. Kohavi, R. A study of cross-validation and bootstrap for accuracy estimation and model selection. In Proceedings of the IJCAI, Montreal, QC, Canada, 20–25 August 1995; Volume 14, pp. 1137–1145.
78. Olofsson, P.; Foody, G.M.; Stehman, S.V.; Woodcock, C.E. Making better use of accuracy data in land change studies: Estimating accuracy and area and quantifying uncertainty using stratified estimation. *Remote Sens. Environ.* **2013**, *129*, 122–131. [[CrossRef](#)]

79. Catal, C. Performance evaluation metrics for software fault prediction studies. *Acta Polytech. Hung.* **2012**, *9*, 193–206.
80. Jeni, L.A.; Cohn, J.F.; De La Torre, F. Facing Imbalanced Data—Recommendations for the Use of Performance Metrics. In Proceedings of the 2013 Humaine Association Conference on Affective Computing and Intelligent Interaction, Geneva, Switzerland, 2–5 September 2013; pp. 245–251.
81. Kuzera, K.; Pontius, R.G. Importance of Matrix Construction for Multiple-Resolution Categorical Map Comparison. *GISci. Remote Sens.* **2008**, *45*, 249–274. [[CrossRef](#)]
82. Pontius, R.G.; Thontteh, O.; Chen, H. Components of information for multiple resolution comparison between maps that share a real variable. *Environ. Ecol. Stat.* **2008**, *15*, 111–142. [[CrossRef](#)]
83. Pontius, R.G.; Parmentier, B. Recommendations for using the relative operating characteristic (ROC). *Landsc. Ecol.* **2014**, *29*, 367–382. [[CrossRef](#)]
84. Pontius, R.G. Component intensities to relate difference by category with difference overall. *Int. J. Appl. Earth Obs. Geoinf.* **2019**, *77*, 94–99. [[CrossRef](#)]
85. Pontius, R.G., Jr.; Connors, J. Range of Categorical Associations for Comparison of Maps with Mixed Pixels. *Photogramm. Eng. Remote Sens.* **2009**, *75*, 963–969. [[CrossRef](#)]
86. Aldwaik, S.Z.; Pontius, R.G., Jr. Intensity analysis to unify measurements of size and stationarity of land changes by interval, category, and transition. *Landsc. Urban Plan.* **2012**, *106*, 103–114. [[CrossRef](#)]
87. Pontius, R.G.; Gao, Y.; Giner, N.M.; Kohyama, T.; Osaki, M.; Hirose, K. Design and Interpretation of Intensity Analysis Illustrated by Land Change in Central Kalimantan, Indonesia. *Land* **2013**, *2*, 351–369. [[CrossRef](#)]
88. Foody, G.M. Harshness in image classification accuracy assessment. *Int. J. Remote Sens.* **2008**, *29*, 3137–3158. [[CrossRef](#)]
89. Cohen, J. A Coefficient of Agreement for Nominal Scales. *Educ. Psychol. Meas.* **1960**, *20*, 37–46. [[CrossRef](#)]
90. Allouche, O.; Tsoar, A.; Kadmon, R. Assessing the Accuracy of Species Distribution Models: Prevalence, Kappa and the True Skill Statistic (TSS). *J. Appl. Ecol.* **2006**, *43*, 1223–1232. [[CrossRef](#)]
91. Foody, G.M. Explaining the unsuitability of the kappa coefficient in the assessment and comparison of the accuracy of thematic maps obtained by image classification. *Remote Sens. Environ.* **2020**, *239*, 111630. [[CrossRef](#)]
92. Willmott, C.J.; Matsuura, K. On the use of dimensioned measures of error to evaluate the performance of spatial interpolators. *Int. J. Geogr. Inf. Sci.* **2006**, *20*, 89–102. [[CrossRef](#)]
93. Willmott, C.J.; Matsuura, K.; Robeson, S.M. Ambiguities inherent in sums-of-squares-based error statistics. *Atmos. Environ.* **2009**, *43*, 749–752. [[CrossRef](#)]
94. Willmott, C.J.; Matsuura, K. Advantages of the mean absolute error (MAE) over the root mean square error (RMSE) in assessing average model performance. *Clim. Res.* **2005**, *30*, 79–82. [[CrossRef](#)]
95. Pontius, R.G., Jr.; Si, K. The total operating characteristic to measure diagnostic ability for multiple thresholds. *Int. J. Geogr. Inf. Sci.* **2014**, *28*, 570–583. [[CrossRef](#)]
96. Fielding, A.H.; Bell, J.F. A review of methods for the assessment of prediction errors in conservation presence/absence models. *Environ. Conserv.* **1997**, *24*, 38–49. [[CrossRef](#)]
97. Blaschke, T. Object based image analysis for remote sensing. *ISPRS J. Photogramm. Remote Sens.* **2010**, *65*, 2–16. [[CrossRef](#)]
98. Costa, H.; Foody, G.M.; Boyd, D.S. Supervised methods of image segmentation accuracy assessment in land cover mapping. *Remote Sens. Environ.* **2018**, *205*, 338–351. [[CrossRef](#)]
99. Powell, R.L.; Matzke, N.; de Souza, C.; Clark, M.; Numata, I.; Hess, L.L.; Roberts, D.A. Sources of error in accuracy assessment of thematic land-cover maps in the Brazilian Amazon. *Remote Sens. Environ.* **2004**, *90*, 221–234. [[CrossRef](#)]
100. Zhong, B.; Ma, P.; Nie, A.; Yang, A.; Yao, Y.; Lü, W.; Zhang, H.; Liu, Q. Land cover mapping using time series HJ-1/CCD data. *Sci. China Earth Sci.* **2014**, *57*, 1790–1799. [[CrossRef](#)]
101. Pacifici, F.; Chini, M.; Emery, W.J. A neural network approach using multi-scale textural metrics from very high-resolution panchromatic imagery for urban land-use classification. *Remote Sens. Environ.* **2009**, *113*, 1276–1292. [[CrossRef](#)]
102. Abbas, I.I.; Muazu, K.M.; Ukoje, J.A. Mapping land use-land cover and change detection in Kafur local government, Katsina, Nigeria (1995-2008) using remote sensing and GIS. *Res. J. Environ. Earth Sci.* **2010**, *2*, 6–12.

103. Sano, E.E.; Rosa, R.; Brito, J.L.S.; Ferreira, L.G. Land cover mapping of the tropical savanna region in Brazil. *Environ. Monit. Assess.* **2010**, *166*, 113–124. [[CrossRef](#)]
104. Hu, T.; Yang, J.; Li, X.; Gong, P. Mapping Urban Land Use by Using Landsat Images and Open Social Data. *Remote Sens.* **2016**, *8*, 151. [[CrossRef](#)]
105. Galletti, C.S.; Myint, S.W. Land-Use Mapping in a Mixed Urban-Agricultural Arid Landscape Using Object-Based Image Analysis: A Case Study from Maricopa, Arizona. *Remote Sens.* **2014**, *6*, 6089–6110. [[CrossRef](#)]
106. Hu, Q.; Wu, W.; Xia, T.; Yu, Q.; Yang, P.; Li, Z.; Song, Q. Exploring the Use of Google Earth Imagery and Object-Based Methods in Land Use/Cover Mapping. *Remote Sens.* **2013**, *5*, 6026–6042. [[CrossRef](#)]
107. Al-Bakri, J.T.; Ajlouni, M.; Abu-Zanat, M. Incorporating Land Use Mapping and Participation in Jordan: An Approach to Sustainable Management of Two Mountainous Areas. *Mt. Res. Dev.* **2008**, *28*, 49–57. [[CrossRef](#)]
108. Liu, J.; Kuang, W.; Zhang, Z.; Xu, X.; Qin, Y.; Ning, J.; Zhou, W.; Zhang, S.; Li, R.; Yan, C.; et al. Spatiotemporal characteristics, patterns, and causes of land-use changes in China since the late 1980s. *J. Geogr. Sci.* **2014**, *24*, 195–210. [[CrossRef](#)]
109. Yadav, P.K.; Kapoor, M.; Sarma, K. Land Use Land Cover Mapping, Change Detection and Conflict Analysis of Nagzira-Navegaon Corridor, Central India Using Geospatial Technology. *Int. J. Remote Sens. GIS* **2012**, *1*.
110. da Costa Freitas, C.; de Souza Soler, L.; Sant'Anna, S.J.S.; Dutra, L.V.; dos Santos, J.R.; Mura, J.C.; Correia, A.H. Land Use and Land Cover Mapping in the Brazilian Amazon Using Polarimetric Airborne P-Band SAR Data. *IEEE Trans. Geosci. Remote Sens.* **2008**, *46*, 2956–2970. [[CrossRef](#)]
111. Dewan, A.M.; Yamaguchi, Y. Land use and land cover change in Greater Dhaka, Bangladesh: Using remote sensing to promote sustainable urbanization. *Appl. Geogr.* **2009**, *29*, 390–401. [[CrossRef](#)]
112. Castañeda, C.; Ducrot, D. Land cover mapping of wetland areas in an agricultural landscape using SAR and Landsat imagery. *J. Environ. Manag.* **2009**, *90*, 2270–2277. [[CrossRef](#)] [[PubMed](#)]
113. Griffiths, P.; van der Linden, S.; Kuemmerle, T.; Hostert, P. A Pixel-Based Landsat Compositing Algorithm for Large Area Land Cover Mapping. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.* **2013**, *6*, 2088–2101. [[CrossRef](#)]
114. Ge, Y. Sub-pixel land-cover mapping with improved fraction images upon multiple-point simulation. *Int. J. Appl. Earth Obs. Geoinf.* **2013**, *22*, 115–126. [[CrossRef](#)]
115. Gong, P.; Wang, J.; Yu, L.; Zhao, Y.; Zhao, Y.; Liang, L.; Niu, Z.; Huang, X.; Fu, H.; Liu, S.; et al. Finer resolution observation and monitoring of global land cover: First mapping results with Landsat TM and ETM+ data. *Int. J. Remote Sens.* **2013**, *34*, 2607–2654. [[CrossRef](#)]
116. Ghorbani, A.; Pakravan, M. Land use mapping using visual vs. digital image interpretation of TM and Google earth derived imagery in Shrivani-Darasi watershed (Northwest of Iran). *Eur. J. Exp. Biol.* **2013**, *3*, 576–582.
117. Deng, J.S.; Wang, K.; Hong, Y.; Qi, J.G. Spatio-temporal dynamics and evolution of land use change and landscape pattern in response to rapid urbanization. *Landsc. Urban Plan.* **2009**, *92*, 187–198. [[CrossRef](#)]
118. Otukey, J.R.; Blaschke, T. Land cover change assessment using decision trees, support vector machines and maximum likelihood classification algorithms. *Int. J. Appl. Earth Obs. Geoinf.* **2010**, *12*, S27–S31. [[CrossRef](#)]
119. Malinverni, E.S.; Tassetti, A.N.; Mancini, A.; Zingaretti, P.; Frontoni, E.; Bernardini, A. Hybrid object-based approach for land use/land cover mapping using high spatial resolution imagery. *Int. J. Geogr. Inf. Sci.* **2011**, *25*, 1025–1043. [[CrossRef](#)]
120. Rozenstein, O.; Karnieli, A. Comparison of methods for land-use classification incorporating remote sensing and GIS inputs. *Appl. Geogr.* **2011**, *31*, 533–544. [[CrossRef](#)]
121. Ran, Y.H.; Li, X.; Lu, L.; Li, Z.Y. Large-scale land cover mapping with the integration of multi-source information based on the Dempster-Shafer theory. *Int. J. Geogr. Inf. Sci.* **2012**, *26*, 169–191. [[CrossRef](#)]
122. Clark, M.L.; Aide, T.M.; Grau, H.R.; Riner, G. A scalable approach to mapping annual land cover at 250 m using MODIS time series data: A case study in the Dry Chaco ecoregion of South America. *Remote Sens. Environ.* **2010**, *114*, 2816–2832. [[CrossRef](#)]
123. Berberoglu, S.; Akin, A. Assessing different remote sensing techniques to detect land use/cover changes in the eastern Mediterranean. *Int. J. Appl. Earth Obs. Geoinf.* **2009**, *11*, 46–53. [[CrossRef](#)]
124. Breiman, L. Random Forests. *Mach. Learn.* **2001**, *45*, 5–32. [[CrossRef](#)]

125. Freeman, E.A.; Moisen, G.G.; Frescino, T.S. Evaluating effectiveness of down-sampling for stratified designs and unbalanced prevalence in Random Forest models of tree species distributions in Nevada. *Ecol. Modell.* **2012**, *233*, 1–10. [[CrossRef](#)]
126. Townshend, J.R.; Masek, J.G.; Huang, C.; Vermote, E.F.; Gao, F.; Channan, S.; Sexton, J.O.; Feng, M.; Narasimhan, R.; Kim, D.; et al. Global characterization and monitoring of forest cover using Landsat data: Opportunities and challenges. *Int. J. Digit. Earth* **2012**, *5*, 373–397. [[CrossRef](#)]
127. Shao, Y.; Lunetta, R.S. Comparison of support vector machine, neural network, and CART algorithms for the land-cover classification using limited training data points. *ISPRS J. Photogramm. Remote Sens.* **2012**, *70*, 78–87. [[CrossRef](#)]
128. Planet Team Planet Application Program Interface. In *Space for Life on Earth*; Planet Labs, Inc.: San Francisco, CA, USA, 2017.
129. Manfreda, S.; McCabe, M.F.; Miller, P.E.; Lucas, R.; Pajuelo Madrigal, V.; Mallinis, G.; Ben Dor, E.; Helman, D.; Estes, L.; Ciraolo, G.; et al. On the Use of Unmanned Aerial Systems for Environmental Monitoring. *Remote Sens.* **2018**, *10*, 641. [[CrossRef](#)]
130. Toutin, T. Geometric processing of IKONOS Geo images with DEM. In Proceedings of the ISPRS Joint Workshop “High Resolution Mapping from Space” 2001, Hanover, Germany, 19–21 September 2001; pp. 19–21.
131. Reinartz, P.; Müller, R.; Schwind, P.; Suri, S.; Bamler, R. Orthorectification of VHR optical satellite data exploiting the geometric accuracy of TerraSAR-X data. *ISPRS J. Photogramm. Remote Sens.* **2011**, *66*, 124–132. [[CrossRef](#)]
132. Aguilar, M.A.; del M. Saldaña, M.; Aguilar, F.J. Assessing geometric accuracy of the orthorectification process from GeoEye-1 and WorldView-2 panchromatic images. *Int. J. Appl. Earth Obs. Geoinf.* **2013**, *21*, 427–435. [[CrossRef](#)]
133. Chen, J.; Zipf, A. DeepVGI: Deep Learning with Volunteered Geographic Information. In *Proceedings of the Proceedings of the 26th International Conference on World Wide Web Companion*; International World Wide Web Conferences Steering Committee: Geneva, Switzerland, 2017; pp. 771–772.
134. Kaiser, P.; Wegner, J.D.; Lucchi, A.; Jaggi, M.; Hofmann, T.; Schindler, K. Learning Aerial Image Segmentation From Online Maps. *IEEE Trans. Geosci. Remote Sens.* **2017**, *55*, 6054–6068. [[CrossRef](#)]
135. Audebert, N.; Le Saux, B.; Lefèvre, S. Joint learning from earth observation and openstreetmap data to get faster better semantic maps. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, Honolulu, HI, USA, 21–26 July 2017; pp. 67–75.
136. Strahler, A.H.; Woodcock, C.E.; Smith, J.A. On the nature of models in remote sensing. *Remote Sens. Environ.* **1986**, *20*, 121–139. [[CrossRef](#)]
137. Foody, G.M. Relating the land-cover composition of mixed pixels to artificial neural network classification output. *Photogramm. Eng. Remote Sens.* **1996**, *62*, 491–498.
138. Moody, A.; Gopal, S.; Strahler, A.H. Artificial neural network response to mixed pixels in coarse-resolution satellite data. *Remote Sens. Environ.* **1996**, *58*, 329–343. [[CrossRef](#)]
139. De Fries, R.S.; Hansen, M.; Townshend, J.R.G.; Sohlberg, R. Global land cover classifications at 8 km spatial resolution: The use of training data derived from Landsat imagery in decision tree classifiers. *Int. J. Remote Sens.* **1998**, *19*, 3141–3168. [[CrossRef](#)]
140. Hansen, M.C.; Potapov, P.V.; Moore, R.; Hancher, M.; Turubanova, S.A.; Tyukavina, A.; Thau, D.; Stehman, S.V.; Goetz, S.J.; Loveland, T.R.; et al. High-resolution global maps of 21st-century forest cover change. *Science* **2013**, *342*, 850–853. [[CrossRef](#)]
141. Kennedy, R.E.; Yang, Z.; Cohen, W.B. Detecting trends in forest disturbance and recovery using yearly Landsat time series: 1. LandTrendr—Temporal segmentation algorithms. *Remote Sens. Environ.* **2010**, *114*, 2897–2910. [[CrossRef](#)]
142. Oppenshaw, S.; Taylor, P. A million or so correlation coefficients. In *Statistical Methods in the Spatial Sciences*; Pion: London, UK, 1979.
143. Jelinski, D.E.; Wu, J. The modifiable areal unit problem and implications for landscape ecology. *Landsc. Ecol.* **1996**, *11*, 129–140. [[CrossRef](#)]
144. Weiss, M.; de Beaufort, L.; Baret, F.; Allard, D.; Bruguier, N.; Marloie, O. Mapping leaf area index measurements at different scales for the validation of large swath satellite sensors: First results of the VALERI project. In Proceedings of the 8th International Symposium in Physical Measurements and Remote Sensing, Aussois, France, 8–12 January 2001; pp. 125–130.

145. Tian, Y.; Woodcock, C.E.; Wang, Y.; Privette, J.L.; Shabanov, N.V.; Zhou, L.; Zhang, Y.; Buermann, W.; Dong, J.; Veikkanen, B.; et al. Multiscale analysis and validation of the MODIS LAI product: I. Uncertainty assessment. *Remote Sens. Environ.* **2002**, *83*, 414–430. [[CrossRef](#)]
146. Masuoka, E.; Roy, D.; Wolfe, R.; Morisette, J.; Sinno, S.; Teague, M.; Saleous, N.; Devadiga, S.; Justice, C.O.; Nickeson, J. MODIS Land Data Products: Generation, Quality Assurance and Validation. In *Land Remote Sensing and Global Environmental Change: NASA's Earth Observing System and the Science of ASTER and MODIS*; Ramachandran, B., Justice, C.O., Abrams, M.J., Eds.; Springer New York: New York, NY, USA, 2011; pp. 509–531. ISBN 9781441967497.
147. Cohen, W.B.; Justice, C.O. Validating MODIS terrestrial ecology products: linking in situ and satellite measurements. *Remote Sens. Environ.* **1999**, *70*, 1–3. [[CrossRef](#)]
148. Fritz, S.; See, L.; McCallum, I.; You, L.; Bun, A.; Moltchanova, E.; Duerauer, M.; Albrecht, F.; Schill, C.; Perger, C.; et al. Mapping global cropland and field size. *Glob. Chang. Biol.* **2015**, *21*, 1980–1992. [[CrossRef](#)] [[PubMed](#)]
149. Debats, S.R.; Estes, L.D.; Thompson, D.R.; Caylor, K.K. Integrating Active Learning and Crowdsourcing into Large-Scale Supervised Landcover Mapping Algorithms. *PeerJ* **2017**. preprints.
150. Estes, L.D.; McRitchie, D.; Choi, J.; Debats, S.; Evans, T.; Guthe, W.; Luo, D.; Ragazzo, G.; Zempleni, R.; Caylor, K.K. A Platform for Crowdsourcing the Creation of Representative, Accurate Landcover Maps. *Environ. Model. Softw.* **2016**, *80*, 41–53. [[CrossRef](#)]
151. Waldner, F.; Schucknecht, A.; Lesiv, M.; Gallego, J.; See, L.; Pérez-Hoyos, A.; d'Andrimont, R.; de Maet, T.; Bayas, J.C.L.; Fritz, S.; et al. Conflation of expert and crowd reference data to validate global binary thematic maps. *Remote Sens. Environ.* **2019**, *221*, 235–246. [[CrossRef](#)]
152. Bey, A.; Sánchez-Paus Díaz, A.; Maniatis, D.; Marchi, G.; Mollicone, D.; Ricci, S.; Bastin, J.-F.; Moore, R.; Federici, S.; Rezende, M.; et al. Collect Earth: Land Use and Land Cover Assessment through Augmented Visual Interpretation. *Remote Sens.* **2016**, *8*, 807. [[CrossRef](#)]
153. Fritz, S.; Sturn, T.; Karner, M.; Moorthy, I.; See, L.; Laso Bayas, J.C.; Fraisl, D. FotoQuest Go: A Citizen Science Approach to the Collection of In-Situ Land Cover and Land Use Data for Calibration and Validation. In Proceedings of the Digital Earth Observation, Salzburg, Austria, 1–4 July 2019.
154. Tuia, D.; Pasolli, E.; Emery, W.J. Using active learning to adapt remote sensing image classifiers. *Remote Sens. Environ.* **2011**, *115*, 2232–2242. [[CrossRef](#)]
155. Van Coillie, F.M.B.; Gardin, S.; Anseel, F.; Duyck, W.; Verbeke, L.P.C.; De Wulf, R.R. Variability of operator performance in remote-sensing image interpretation: the importance of human and external factors. *Int. J. Remote Sens.* **2014**, *35*, 754–778. [[CrossRef](#)]
156. Johnson, B.A.; Iizuka, K. Integrating OpenStreetMap crowdsourced data and Landsat time-series imagery for rapid land use/land cover (LULC) mapping: Case study of the Laguna de Bay area of the Philippines. *Appl. Geogr.* **2016**, *67*, 140–149. [[CrossRef](#)]
157. Neigh, C.S.R.; Carroll, M.L.; Wooten, M.R.; McCarty, J.L.; Powell, B.F.; Husak, G.J.; Enekel, M.; Hain, C.R. Smallholder crop area mapped with wall-to-wall WorldView sub-meter panchromatic image texture: A test case for Tigray, Ethiopia. *Remote Sens. Environ.* **2018**, *212*, 8–20. [[CrossRef](#)]
158. Clark, M.L.; Aide, T.M.; Riner, G. Land change for all municipalities in Latin America and the Caribbean assessed from 250-m MODIS imagery (2001–2010). *Remote Sens. Environ.* **2012**, *126*, 84–103. [[CrossRef](#)]
159. Comber, A.; Fisher, P.; Wadsworth, R. What is land cover? *Environ. Plan.* **2005**, *32*, 199–209. [[CrossRef](#)]
160. Kohli, D.; Sliuzas, R.; Kerle, N.; Stein, A. An ontology of slums for image-based classification. *Comput. Environ. Urban Syst.* **2012**, *36*, 154–163. [[CrossRef](#)]
161. Verburg, P.H.; Neumann, K.; Nol, L. Challenges in using land use and land cover data for global change studies. *Glob. Chang. Biol.* **2011**, *17*, 974–989. [[CrossRef](#)]
162. Weng, Q. Remote sensing of impervious surfaces in the urban areas: Requirements, methods, and trends. *Remote Sens. Environ.* **2012**, *117*, 34–49. [[CrossRef](#)]
163. Kohli, D.; Stein, A.; Sliuzas, R. Uncertainty analysis for image interpretations of urban slums. *Comput. Environ. Urban Syst.* **2016**, *60*, 37–49. [[CrossRef](#)]
164. Rocchini, D. While Boolean sets non-gently rip: A theoretical framework on fuzzy sets for mapping landscape patterns. *Ecol. Complex.* **2010**, *7*, 125–129. [[CrossRef](#)]
165. Woodcock, C.E.; Gopal, S. Fuzzy set theory and thematic maps: Accuracy assessment and area estimation. *Int. J. Geogr. Inf. Sci.* **2000**, *14*, 153–172. [[CrossRef](#)]

166. Rocchini, D.; Foody, G.M.; Nagendra, H.; Ricotta, C.; Anand, M.; He, K.S.; Amici, V.; Kleinschmit, B.; Förster, M.; Schmidtlein, S.; et al. Uncertainty in ecosystem mapping by remote sensing. *Comput. Geosci.* **2013**, *50*, 128–135. [[CrossRef](#)]
167. Zhang, J.; Foody, G.M. A fuzzy classification of sub-urban land cover from remotely sensed imagery. *Int. J. Remote Sens.* **1998**, *19*, 2721–2738. [[CrossRef](#)]
168. Woodcock, C.E.; Strahler, A.H. The factor of scale in remote sensing. *Remote Sens. Environ.* **1987**, *21*, 311–332. [[CrossRef](#)]
169. Cracknell, A.P. Review article Synergy in remote sensing—what’s in a pixel? *Int. J. Remote Sens.* **1998**, *19*, 2025–2047. [[CrossRef](#)]
170. Pontius, R.G.; Cheuk, M.L. A generalized cross-tabulation matrix to compare soft-classified maps at multiple resolutions. *Int. J. Geogr. Inf. Sci.* **2006**, *20*, 1–30. [[CrossRef](#)]
171. Silván-Cárdenas, J.L.; Wang, L. Sub-pixel confusion–uncertainty matrix for assessing soft classifications. *Remote Sens. Environ.* **2008**, *112*, 1081–1095. [[CrossRef](#)]
172. Foody, G.M. The continuum of classification fuzziness in thematic mapping. *Photogramm. Eng. Remote Sens.* **1999**, *65*, 443–452.
173. Foody, G.M. Fully fuzzy supervised classification of land cover from remotely sensed imagery with an artificial neural network. *Neural Comput. Appl.* **1997**, *5*, 238–247. [[CrossRef](#)]
174. Laso Bayas, J.C.; See, L.; Fritz, S.; Sturn, T.; Perger, C.; Dürauer, M.; Karner, M.; Moorthy, I.; Schepaschenko, D.; Domian, D.; et al. Crowdsourcing In-Situ Data on Land Cover and Land Use Using Gamification and Mobile Technology. *Remote Sens.* **2016**, *8*, 905. [[CrossRef](#)]
175. Tewkesbury, A.P.; Comber, A.J.; Tate, N.J.; Lamb, A.; Fisher, P.F. A critical synthesis of remotely sensed optical image change detection techniques. *Remote Sens. Environ.* **2015**, *160*, 1–14. [[CrossRef](#)]
176. Stehman, S.V.; Fonte, C.C.; Foody, G.M.; See, L. Using volunteered geographic information (VGI) in design-based statistical inference for area estimation and accuracy assessment of land cover. *Remote Sens. Environ.* **2018**, *212*, 47–59. [[CrossRef](#)]
177. Thompson, I.D.; Maher, S.C.; Rouillard, D.P.; Fryxell, J.M.; Baker, J.A. Accuracy of forest inventory mapping: Some implications for boreal forest management. *For. Ecol. Manag.* **2007**, *252*, 208–221. [[CrossRef](#)]
178. Bland, M.J.; Altman, D.G. Statistics notes: Measurement error. *BMJ* **1996**, *312*, 1654. [[CrossRef](#)] [[PubMed](#)]
179. Martin, D. An Introduction to “The Guide to the Expression of Uncertainty in Measurement”. In *Evaluation of Measurement Data—Guide to the Expression of Uncertainty in Measurement*; JCGM: Geneva, Switzerland, 2008; pp. 1–10.
180. Thanh Noi, P.; Kappas, M. Comparison of Random Forest, k-Nearest Neighbor, and Support Vector Machine Classifiers for Land Cover Classification Using Sentinel-2 Imagery. *Sensors* **2017**, *18*. [[CrossRef](#)] [[PubMed](#)]
181. Song, K. Tackling Uncertainties and Errors in the Satellite Monitoring of Forest Cover Change. Ph.D. Thesis, University of Maryland, College Park, MD, USA, 2010.
182. Foody, G.M. The impact of imperfect ground reference data on the accuracy of land cover change estimation. *Int. J. Remote Sens.* **2009**, *30*, 3275–3281. [[CrossRef](#)]
183. Foody, G.M. Ground reference data error and the mis-estimation of the area of land cover change as a function of its abundance. *Remote Sens. Lett.* **2013**, *4*, 783–792. [[CrossRef](#)]
184. Homer, C.G.; Fry, J.A.; Barnes, C.A.; National land cover dataset (NLCD). *The National Land Cover Database*; U.S. Geological Survey: Reston, VA, USA, 2012.
185. Menon, S.; Akbari, H.; Mahanama, S.; Sednev, I.; Levinson, R. Radiative forcing and temperature response to changes in urban albedos and associated CO<sub>2</sub> offsets. *Environ. Res. Lett.* **2010**, *5*, 014005. [[CrossRef](#)]
186. Hutyrá, L.R.; Yoon, B.; Hepinstall-Cymerman, J.; Alberti, M. Carbon consequences of land cover change and expansion of urban lands: A case study in the Seattle metropolitan region. *Landsc. Urban Plan.* **2011**, *103*, 83–93. [[CrossRef](#)]
187. Reinmann, A.B.; Hutyrá, L.R.; Trlica, A.; Olofsson, P. Assessing the global warming potential of human settlement expansion in a mesic temperate landscape from 2005 to 2050. *Sci. Total Environ.* **2016**, *545–546*, 512–524. [[CrossRef](#)]
188. Hardiman, B.S.; Wang, J.A.; Hutyrá, L.R.; Gately, C.K.; Getson, J.M.; Friedl, M.A. Accounting for urban biogenic fluxes in regional carbon budgets. *Sci. Total Environ.* **2017**, *592*, 366–372. [[CrossRef](#)]
189. Seto, K.C.; Güneralp, B.; Hutyrá, L.R. Global forecasts of urban expansion to 2030 and direct impacts on biodiversity and carbon pools. *Proc. Natl. Acad. Sci. USA* **2012**, *109*, 16083–16088. [[CrossRef](#)]

190. Angel, S.; Parent, J.; Civco, D.L.; Blei, A.; Potere, D. The dimensions of global urban expansion: Estimates and projections for all countries, 2000–2050. *Prog. Plann.* **2011**, *75*, 53–107. [[CrossRef](#)]
191. Coulston, J.W.; Moisen, G.G.; Wilson, B.T.; Finco, M.V.; Cohen, W.B.; Brewer, C.K. Modeling percent tree canopy cover: A pilot study. *Photogramm. Eng. Remote Sens.* **2012**, *78*, 715–727. [[CrossRef](#)]
192. Reinmann, A.B.; Hutyra, L.R. Edge effects enhance carbon uptake and its vulnerability to climate change in temperate broadleaf forests. *Proc. Natl. Acad. Sci. USA* **2017**, *114*, 107–112. [[CrossRef](#)] [[PubMed](#)]
193. Rolnick, D.; Veit, A.; Belongie, S.; Shavit, N. Deep Learning is Robust to Massive Label Noise. *arXiv* **2017**.
194. Nachmany, Y.; Alemohammad, H. Detecting Roads from Satellite Imagery in the Developing World. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, Long Beach, CA, USA, 16–20 June 2019; pp. 83–89.
195. The SpaceNet Catalog SpaceNet on Amazon Web Services (AWS). “Datasets.” The SpaceNet Catalog. Last modified 30 April 2018. Available online: <https://spacenetchallenge.github.io/datasets/datasetHomePage.html> (accessed on 15 November 2019).
196. Alemohammad, S.H.; Fang, B.; Konings, A.G.; Aires, F.; Green, J.K.; Kolassa, J.; Miralles, D.; Prigent, C.; Gentine, P. Water, Energy, and Carbon with Artificial Neural Networks (WECANN): A statistically-based estimate of global surface turbulent fluxes and gross primary productivity using solar-induced fluorescence. *Biogeosciences* **2017**, *14*, 4101–4124. [[CrossRef](#)] [[PubMed](#)]
197. McColl, K.A.; Vogelzang, J.; Konings, A.G.; Entekhabi, D.; Piles, M.; Stoffelen, A. Extended triple collocation: Estimating errors and correlation coefficients with respect to an unknown target. *Geophys. Res. Lett.* **2014**, *41*, 6229–6236. [[CrossRef](#)]
198. Debats, S.R.; Luo, D.; Estes, L.D.; Fuchs, T.J.; Caylor, K.K. A Generalized Computer Vision Approach to Mapping Crop Fields in Heterogeneous Agricultural Landscapes. *Remote Sens. Environ.* **2016**, *179*, 210–221. [[CrossRef](#)]
199. Estes, L.D.; Ye, S.; Song, L.; Avery, R.B.; McRitchie, D.; Eastman, J.R.; Debats, S.R. *Improving Maps of Smallholder-Dominated Croplands through Tight Integration of Human and Machine Intelligence*; American Geophysical Union: Washington, DC, USA, 2019.
200. Jain, M.; Balwinder-Singh; Rao, P.; Srivastava, A.K.; Poonia, S.; Blesh, J.; Azzari, G.; McDonald, A.J.; Lobell, D.B. The impact of agricultural interventions can be doubled by using satellite data. *Nat. Sustain.* **2019**, *2*, 931–934. [[CrossRef](#)]
201. Pontius, R.G. Criteria to Confirm Models that Simulate Deforestation and Carbon Disturbance. *Land* **2018**, *7*, 105. [[CrossRef](#)]
202. Schennach, S.M. Recent Advances in the Measurement Error Literature. *Annu. Rev. Econom.* **2016**, *8*, 341–377. [[CrossRef](#)]
203. Waldner, F.; De Aballeyra, D.; Verón, S.R.; Zhang, M.; Wu, B.; Plotnikov, D.; Bartalev, S.; Lavreniuk, M.; Skakun, S.; Kussul, N.; et al. Towards a set of agrosystem-specific cropland mapping methods to address the global cropland diversity. *Int. J. Remote Sens.* **2016**, *37*, 3196–3231. [[CrossRef](#)]
204. Castelluccio, M.; Poggi, G.; Sansone, C.; Verdoliva, L. Land Use Classification in Remote Sensing Images by Convolutional Neural Networks. *arXiv* **2015**.
205. Azevedo, T., Sr.; Souza, C.M., Jr.; Shimbo, J.; Alencar, A. *MapBiomass Initiative: Mapping Annual Land Cover and Land Use Changes in Brazil from 1985 to 2017*; American Geophysical Union: Washington, DC, USA, 2018; Volume 2018.
206. Brown, J.F.; Tollerud, H.J.; Barber, C.P.; Zhou, Q.; Dwyer, J.L.; Vogelmann, J.E.; Loveland, T.R.; Woodcock, C.E.; Stehman, S.V.; Zhu, Z.; et al. Lessons learned implementing an operational continuous United States national land change monitoring capability: The Land Change Monitoring, Assessment, and Projection (LCMAP) approach. *Remote Sens. Environ.* **2019**, 111356. [[CrossRef](#)]
207. Estes, L.; Elsen, P.R.; Treuer, T.; Ahmed, L.; Caylor, K.; Chang, J.; Choi, J.J.; Ellis, E.C. The spatial and temporal domains of modern ecology. *Nat. Ecol. Evol.* **2018**, *2*, 819–826. [[CrossRef](#)] [[PubMed](#)]
208. Jensen, J.R.; Cowen, D.C. Remote sensing of urban/suburban infrastructure and socio-economic attributes. *Photogramm. Eng. Remote Sens.* **1999**, *65*, 611–622.
209. Dorais, A.; Cardille, J. Strategies for Incorporating High-Resolution Google Earth Databases to Guide and Validate Classifications: Understanding Deforestation in Borneo. *Remote Sens.* **2011**, *3*, 1157–1176. [[CrossRef](#)]
210. Sexton, J.O.; Urban, D.L.; Donohue, M.J.; Song, C. Long-term land cover dynamics by multi-temporal classification across the Landsat-5 record. *Remote Sens. Environ.* **2013**, *128*, 246–258. [[CrossRef](#)]



211. Reis, M.S.; Escada, M.I.S.; Dutra, L.V.; Sant'Anna, S.J.S.; Vogt, N.D. Towards a Reproducible LULC Hierarchical Class Legend for Use in the Southwest of Pará State, Brazil: A Comparison with Remote Sensing Data-Driven Hierarchies. *Land* **2018**, *7*, 65. [[CrossRef](#)]
212. Anderson, J.R. *A Land Use and Land Cover Classification System for Use with Remote Sensor Data*; U.S. Government Printing Office: Washington, DC, USA, 1976.
213. Herold, M.; Woodcock, C.E.; di Gregorio, A.; Mayaux, P.; Belward, A.S.; Latham, J.; Schullius, C.C. A joint initiative for harmonization and validation of land cover datasets. *IEEE Trans. Geosci. Remote Sens.* **2006**, *44*, 1719–1727. [[CrossRef](#)]
214. Carletto, C.; Gourlay, S.; Winters, P. From Guesstimates to GPStimates: Land Area Measurement and Implications for Agricultural Analysis. *J. Afr. Econ.* **2015**, *24*, 593–628. [[CrossRef](#)]
215. See, L.; Comber, A.; Salk, C.; Fritz, S.; van der Velde, M.; Perger, C.; Schill, C.; McCallum, I.; Kraxner, F.; Obersteiner, M. Comparing the quality of crowdsourced data contributed by expert and non-experts. *PLoS ONE* **2013**, *8*, e69958. [[CrossRef](#)]
216. Phinn, S.R. A framework for selecting appropriate remotely sensed data dimensions for environmental monitoring and management. *Int. J. Remote Sens.* **1998**, *19*, 3457–3463. [[CrossRef](#)]
217. Phinn, S.R.; Menges, C.; Hill, G.J.E.; Stanford, M. Optimizing Remotely Sensed Solutions for Monitoring, Modeling, and Managing Coastal Environments. *Remote Sens. Environ.* **2000**, *73*, 117–132. [[CrossRef](#)]
218. Lu, D.; Weng, Q. A survey of image classification methods and techniques for improving classification performance. *Int. J. Remote Sens.* **2007**, *28*, 823–870. [[CrossRef](#)]
219. Cingolani, A.M.; Renison, D.; Zak, M.R.; Cabido, M.R. Mapping vegetation in a heterogeneous mountain rangeland using landsat data: An alternative method to define and classify land-cover units. *Remote Sens. Environ.* **2004**, *92*, 84–97. [[CrossRef](#)]
220. Burke, M.; Lobell, D.B. Satellite-based assessment of yield variation and its determinants in smallholder African systems. *Proc. Natl. Acad. Sci. USA* **2017**. [[CrossRef](#)] [[PubMed](#)]
221. Jin, Z.; Azzari, G.; You, C.; Di Tommaso, S.; Aston, S.; Burke, M.; Lobell, D.B. Smallholder maize area and yield mapping at national scales with Google Earth Engine. *Remote Sens. Environ.* **2019**, *228*, 115–128. [[CrossRef](#)]
222. Lobell, D.B.; Thau, D.; Seifert, C.; Engle, E.; Little, B. A Scalable Satellite-Based Crop Yield Mapper. *Remote Sens. Environ.* **2015**, *164*, 324–333. [[CrossRef](#)]
223. Grassini, P.; van Bussel, L.G.J.; Van Wart, J.; Wolf, J.; Claessens, L.; Yang, H.; Boogaard, H.; de Groot, H.; van Ittersum, M.K.; Cassman, K.G. How Good Is Good Enough? Data Requirements for Reliable Crop Yield Simulations and Yield-Gap Analysis. *Field Crops Res.* **2015**, *177*, 49–63. [[CrossRef](#)]
224. Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; et al. ImageNet Large Scale Visual Recognition Challenge. *Int. J. Comput. Vis.* **2015**, *115*, 211–252. [[CrossRef](#)]
225. Fu, X.; McCane, B.; Mills, S.; Albert, M. NOKMeans: Non-Orthogonal K-means Hashing. In *Computer Vision—ACCV 2014*; Cremers, D., Reid, I., Saito, H., Yang, M.-H., Eds.; Lecture Notes in Computer Science; Springer International Publishing: Cham, Switzerland, 2015; Volume 9003, pp. 162–177. ISBN 9783319168647.
226. Basu, S.; Ganguly, S.; Mukhopadhyay, S.; DiBiano, R.; Karki, M.; Nemani, R. DeepSat: A Learning Framework for Satellite Imagery. In Proceedings of the Proceedings of the 23rd SIGSPATIAL International Conference on Advances in Geographic Information Systems, Seattle, WA, USA, 3–6 November 2015; ACM: New York, NY, USA, 2015; pp. 37:1–37:10.
227. Yang, Y.; Newsam, S. Bag-of-visual-words and spatial extensions for land-use classification. In Proceedings of the 18th SIGSPATIAL International Conference on Advances in Geographic Information Systems, San Jose, CA, USA, 2 November 2010.
228. Shen, C. A Transdisciplinary Review of Deep Learning Research and Its Relevance for Water Resources Scientists. *Water Resour. Res.* **2018**, *54*, 8558–8593. [[CrossRef](#)]

