

1 **Implementation of CCDC to produce the LCMAP Collection 1.0 annual land surface**
2 **change product**

3 George Z. Xian¹, Kelcy Smith², Danika Wellington², Josephine Horton², Qiang Zhou³,
4 Congcong Li³, Roger Auch¹, Jesslyn F. Brown¹, Zhe Zhu⁴, and Ryan R. Reker²

5 ¹United States Geological Survey (USGS) Earth Resources Observation and Science (EROS)
6 Center, Sioux Falls, South Dakota 57198, U.S.A.

7 ²KBR, Contractor to the USGS EROS Center, Sioux Falls, SD, 57198, U.S.A.

8 ³ASRC Federal Data Solutions (AFDS), Contractor to the USGS EROS Sioux Falls, SD 57198,
9 U.S.A.

10 ⁴Department of Natural Resources and the Environment, University of Connecticut, Storrs, CT,
11 U.S.A.

12 Correspondence: George Xian (xian@usgs.gov)

13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30

31 **Abstract**

32 The increasing availability of high-quality remote sensing data and advanced technologies has
33 spurred land cover mapping to characterize land change from local to global scales. However,
34 most land change datasets either span multiple decades at a local scale or cover limited time over
35 a larger geographic extent. Here, we present a new land cover and land surface change dataset
36 created by the Land Change Monitoring, Assessment, and Projection (LCMAP) program over
37 the conterminous United States (CONUS). The LCMAP land cover change dataset consists of
38 annual land cover and land cover change products over the period 1985-2017 at 30-m resolution
39 using Landsat and other ancillary data via the Continuous Change Detection and Classification
40 (CCDC) algorithm. In this paper, we describe our novel approach to implement the CCDC
41 algorithm to produce the LCMAP product suite composed of five land cover and five land
42 surface change related products. The LCMAP land cover products were validated using a
43 collection of ~ 25,000 reference samples collected independently across CONUS. The overall
44 agreement for all years of the LCMAP primary land cover product reached 82.5%. The LCMAP
45 products are produced through the LCMAP Information Warehouse and Data Store (IW+DS)
46 and Shared Mesos Cluster systems that can process, store, and deliver all datasets for public
47 access. To our knowledge, this is the first set of published 30-m annual land change datasets that
48 include land cover, land cover change, and spectral change spanning from the 1980s to the
49 present for the United States. The LCMAP product suite provides useful information for land
50 resource management and facilitates studies to improve the understanding of terrestrial
51 ecosystems and the complex dynamics of the Earth system. The LCMAP system could be
52 implemented to produce global land change products in the future.

53
54
55
56
57
58
59
60

61 **1 Introduction**

62

63 Changes in land cover and land surface are one of the greatest and most immediate influences on
64 the Earth system, and these changes will continue in association with a surging human
65 population and growing demand on land resources (Szantoi et al., 2020). Changes in land cover
66 and ecosystems and their implications for global environmental change and sustainability are
67 major research challenges for developing strategies to respond to ongoing global change while
68 meeting development goals (Turner II et al., 2007). Unknowns related to the spatial extent and
69 degrees of impacts of anthropogenic activities on natural systems and strategies to respond to
70 ongoing global change hinder efforts to overcome sustainability challenges (Erb et al., 2017;
71 Reid et al., 2010). An improved understanding of the complex and dynamic interactions between
72 the various Earth system components, including humans and their activities, is critical for
73 policymakers and scientists (Foley, 2005; Foley et al., 2011). To fully understand these processes
74 and monitor these changes, accurate and frequently updated land cover information is essential
75 for scientific research and to assist decision makers in responding to the challenges associated
76 with competing land demands and land surface change.

77 The characteristics of land surface fundamentally connect with the functioning of Earth's
78 terrestrial surface. Satellite observations have been used to observe the Earth's surface and to
79 characterize land cover and change from local to global scales. Remote sensing data allows us to
80 obtain information over large areas in a practical and accurate manner. With advanced
81 technologies and accumulating satellite data, countries and regions have produced multi-spatial
82 and multi-temporal resolution land cover products (Chen et al., 2015; Gong et al., 2020; Hansen,
83 2013; Homer et al., 2020; Li et al., 2020). A variety of land change mapping has been carried out
84 to produce land cover and change products in the United States. Among these efforts are the
85 widely known National Land Cover Database (NLCD) products. NLCD has provided
86 comprehensive, general-purpose land cover mapping products at 30-m resolution since 2001 in
87 the United States, and the products have been published and updated across more than a decade
88 (Homer et al., 2020). NLCD provides Anderson Level II land cover classification (Anderson,
89 1976) for the conterminous United States (CONUS) at approximately 2–3-year intervals. Other
90 national-scale mapping projects focus on specific land cover themes. Among these are the
91 Landscape Fire and Resource Management Planning Tools (LANDFIRE) (Picotte et al., 2019),

92 which maps vegetation and fuels in support of wildfire management, and the Cropland Data
93 Layer (Boryan et al., 2011) generated by the National Agricultural Statistics Service (NASS) of
94 the United States Department of Agriculture (USDA). Due to the need to incorporate data from
95 neighboring years, as well as extensive post-processing, ancillary dataset dependencies, and
96 analyst-supported refinement, release dates for both LANDFIRE and NLCD products are
97 typically several years subsequent to the nominal map year. Other products including national
98 urban extent change and vegetation phenology data are available (Li et al., 2019; Li et al., 2020).
99 These projects vary in how land change information is incorporated or expressed across product
100 releases. Continuous data stacks allow for an increase in input features for land cover
101 classification. Frequent data also provides the opportunity for near-real time change monitoring
102 with frequently updated image acquisitions. The availability of land change information has led
103 to approaches that attempt to monitor surface properties continuously through time (Franklin et
104 al., 2015; Gong et al., 2019; Hermosilla et al., 2018; Homer et al., 2020; Kennedy et al., 2015; Li
105 et al., 2020). Such approaches have several advantages over traditional image processing
106 techniques based on small numbers of images (Bullock et al., 2020; Zhu and Woodcock, 2014b).
107 Leveraging the increasingly massive amount of openly available, analysis-ready data products
108 into the generation of operational land cover and land change information has been described as
109 the new paradigm for land cover science (Wulder et al., 2018). The approach, which intended to
110 use all available medium resolution remotely sensed data from the 1980s to the present, opened a
111 door for the scientific community to integrate time series information to improve change
112 detection and land cover characterization in a robust way. Furthermore, change events, when
113 combined with knowledge of ecology settings or anticipation of a given process post-change, can
114 accommodate consistent change observations and characterization of land cover. For example,
115 forest areas that are cleared by wildfire or harvest activities typically transfer to non-forest
116 herbaceous or shrub vegetation cover, followed by a succession of young tree stages, ultimately
117 returning to a forest class. Traditional change detection methods using limited observations may
118 not have identified these changes if data were collected with a starting date prior to the change
119 and an ending date that occurred after the transitional (non-tree) vegetation returned to tree
120 cover. Therefore, incorporating change information into the land cover characterization process
121 allows for insights regarding expected land cover class transitions related to successional
122 processes, and likewise provides a mechanism to identify illogical class transitions and cause or

123 agent of change (Kennedy et al., 2015; Wulder et al., 2018). The choice of a time series
124 approach also allows missing data and phenological variations to be handled robustly (Friedl et
125 al., 2010; Wulder et al., 2018).

126 The Continuous Change Detection (CCD) and Classification (CCDC) algorithm (Zhu and
127 Woodcock, 2014b; Zhu et al., 2015b) was developed to advance time series change detection by
128 using all available Landsat data. The CCD algorithm uses robust methodology to identify when
129 and how the land surface changes through time. The algorithm first estimates a time series model
130 based on clear observations and then detects outliers by comparing model estimates and Landsat
131 observations. The algorithm fits harmonic regression models through a Least Absolute Shrinkage
132 and Selection Operator (LASSO) (Tibshirani, 1996) approach to every pixel over time to
133 estimate the time series model defined by sine and cosine functions. New Landsat records are
134 compared to predicted results, and if the observed data deviate beyond a set threshold for all
135 records within a moving window period, then a model break is produced. The parameters used to
136 fit the model are used as inputs for the cover classifier for land cover characterization.

137 The original implementation of CCDC was written in the MATLAB programming language and
138 had been implemented for a regional land cover change assessment in the eastern CONUS (Zhu
139 and Woodcock, 2014b). The algorithm includes the automation of change detection/classification
140 and can monitor changes for different land cover types. The implementation of CCDC into a
141 large geographic extent still encounters several challenges: the availability of Landsat records
142 and training datasets, the effectiveness of choosing good quality Landsat records, and the
143 robustness to characterize land cover and change across various land cover types and conditions.
144 In this paper, we outlined major efforts and challenges in the implementation of CCDC for the
145 U.S. Geological Survey (USGS) Land Change Monitoring, Assessment, and Projection
146 (LCMAP) initiative (Brown et al., 2020). LCMAP focuses on using CCD/CCDC with time series
147 Landsat records and other ancillary information to produce annual land cover and change
148 products from 1985 to the present for the United States. We focused on how LCMAP employed
149 every observation in a time series of U.S. Landsat Analysis Ready Data (ARD) (Dwyer et al.,
150 2018) over a long period starting with the 1980s to determine whether change occurred at any
151 given point in the observation record. The CCDC algorithm that was initially developed for
152 abrupt change detection on the land surface was modified through lessons learned from the

153 prototype test to include both gradual land cover transition and abrupt land change so that the
154 algorithm could be used in an operational setting with the goals of robust, repeatable, and
155 geographically consistent results (Brown et al., 2020). The algorithm was further used to classify
156 the pixel to indicate what land cover type(s) were observed before and after a detected change on
157 the land surface. Classification in LCMAP was modified to improve representativeness of
158 training data and reduce notable artifacts including misclassification of rare classes and dramatic
159 increase in the amount of training data. The CCDC algorithm has since been translated into an
160 open-source library as Python code. The full implementation joined the CCD Python library with
161 the classification methodology in combination with data delivery/processing services made
162 available through the LCMAP Information Warehouse and Data Store (IW+DS) and evolved as
163 a national operational monitoring system.

164

165 **2 Data Sources**

166 The CCDC algorithm utilizes all available Landsat observations including surface reflectance,
167 brightness temperature, and associated quality data to characterize the spectral responses of
168 every pixel through harmonic regression model fits. The model fits are then used to categorize
169 each pixel time series into temporal segments of stable periods and to estimate the dates at which
170 the spectral time-series data diverge from past responses or patterns. The outcomes of model fits
171 and other input data are then used for classification. The algorithm requires several input datasets
172 to perform both change detection and classification.

173 **2.1 Landsat observations**

174 U.S. Landsat ARD have been processed to a minimum set of requirements and organized into a
175 form that can be more directly and easily used for monitoring and assessing landscape change
176 with minimal additional user effort. Landsat ARD Collection 1 provides consistent radiometric
177 and geometric Landsat products across Landsat 4-5 Thematic Mapper (TM), Landsat 7 Enhanced
178 Thematic Mapper Plus (ETM+), and Landsat 8 Operational Land Imager (OLI) / Thermal
179 Infrared Sensor (TIRS) instruments for use in time series analysis (Dwyer et al., 2018). Landsat
180 ARD is organized in tiles, which are units of uniform dimension bounded by static corner points
181 in a defined grid system (Fig. 1). An ARD tile is currently defined as 5,000 x 5,000 30-meter (m)
182 pixels or 150 x 150-kilometer (km). To implement CCDC algorithms to produce LCMAP

183 Collection 1.0 land change products in CONUS, all available Landsat ARD records of surface
184 reflectance and brightness temperature from the 1980s to 2017 were required.

185 **2.2 Land cover and ancillary datasets**

186 The CCDC algorithm employs every observation in a time series of Landsat data to determine
187 whether change has occurred at any given time. The algorithm further classifies the time series to
188 indicate what land cover types were observed before and after a detected change and further to
189 generate LCMAP annual land cover products (Table 1). The land cover products are produced by
190 using training data from NLCD in 2001. NLCD provides Anderson Level II (Anderson, 1976)
191 land cover classification for CONUS and outlying areas (Homer et al., 2020). Spectral index and
192 change metrics between cloud-corrected Landsat mosaics are used, among other information, to
193 identify change pixels (Jin et al., 2013). These metrics allow NLCD to incorporate temporal and
194 spectral trajectory information into both training data selection and final land cover
195 classification. The NLCD land cover data is used in LCMAP as land cover training data.

196
197 Ancillary data comprises two main source datasets: the USGS National Elevation Dataset (NED)
198 (Gesch et al., 2002) 1 arc-second Digital Elevation Models (DEM), and a wetland potential index
199 (WPI) layer created for NLCD 2011 land cover production (Zhu et al., 2016). The WPI layer is a
200 ranking (0–8) of wetland likelihood from a comparison of the National Wetland Inventory
201 (NWI), the U.S. Department of Agriculture Soil Survey Geographic Database (SSURGO) for
202 hydric soils, and the NLCD 2006 wetlands land cover classes.

204 **3 Methodology**

205 As part of the operational LCMAP system, the original MATLAB version of the CCDC
206 algorithm is converted to a format that meets the needs of large-scale land change detection and
207 change characterization on an annual basis. Python is selected to replace MATLAB to implement
208 the CCDC algorithm for LCMAP. The CCD component of the CCDC algorithm is converted to
209 create the Python-based CCD (PyCCD) library. The PyCCD library is a per-pixel algorithm, and
210 the fundamental outputs are the spectral characterizations (segments) of the input data. There are
211 several key components in PyCCD. The overall CCD procedures are summarized in Fig. 2.

212 **3.1 Data filtering and Harmonic modeling**

213 The removal of invalid and cloud-contaminated data points is important for deriving model
 214 coefficients that accurately represent the phenology of the surface, and for the correct
 215 identification of model break points. The CCD algorithm uses Landsat ARD PIXELQA values to
 216 mask observations identified as cloud, cloud shadow, fill, or (in some cases) snow derived based
 217 on the Fmask 3.3 algorithm (Zhu et al., 2015a; Zhu and Woodcock, 2012). Additional cirrus and
 218 terrain occlusion bits are provided for Landsat 8 OLI-TIRS ARD that are not available in the
 219 Landsat 4–7 TM/ETM+ quality assessment band. To maintain consistency across the historical
 220 archive, the algorithm does not rely on these Landsat 8-only QA flags to filter out observations.

221 Landsat ARD containing invalid or physically unrealistic data values are removed. For the
 222 surface reflectance bands, the valid data range is between 0 and 10000. Brightness temperature
 223 values, which in the ARD are stored as $10 \times$ temperature (kelvin), are converted to $100 \times$ °C and
 224 observations are filtered for values outside the range -9320 and 7070 (-93.2–70.7°C). This
 225 procedure rescales the brightness temperature values into a roughly similar numerical range as
 226 the surface reflectance bands. A multitemporal mask (Tmask) model (Zhu and Woodcock,
 227 2014a) is implemented first to remove additional outliers by using the multitemporal observation
 228 record to identify values that deviate from the overall phenology curve using a specific harmonic
 229 model to perform an initial fit to the phenology. Additional details are provided in the
 230 Supplementary materials S1.

231 The filtered Landsat ARD is further operated to generate the time series fit by harmonic models
 232 whose sinusoidal components are frequency multiples of the base annual frequency. A constant
 233 and linear term characterizes the surface reflectance or brightness temperature offset value and
 234 overall slope, respectively. The full harmonic model is defined as follows:

$$235 \hat{p}(i, t) = c_{0,i} + c_{1,i}t + \sum_{n=1}^3 (a_{n,i} \cos \omega nt + b_{n,i} \sin \omega nt) \quad (1)$$

236 where ω is the base annual frequency ($2\pi/T$), t is the ordinal of the date when January 1 of the
 237 year zero has ordinal 1 (sometimes called Julian date), i is the i th Landsat band, $a_{n,i}$ and $b_{n,i}$ are
 238 the estimated n th order harmonic coefficients for the i th Landsat band, $c_{0,i}$ and $c_{1,i}$ are the
 239 estimated intercept and slope coefficients for the i th Landsat band, and $\hat{p}(i, t)$ is the predicted
 240 value for the i th Landsat band at ordinal date t . Model initialization and certain special-case
 241 regression fits such as at the beginning/end of the time series use the simple four-coefficient

242 model. Outside of these conditions, the selection of coefficient depends on the number of
243 observations used for the regression. For a full model (eight coefficients), there must be at least
244 24 observations covered by the regression. The fit parameters returned by PyCCD always
245 include eight coefficient values including an intercept, with unused coefficients reported as
246 zeroes.

247 **3.2 Regression models and change detection thresholds**

248 The best-fit coefficients for the time series model are calculated using a LASSO regression
249 model (Tibshirani, 1996). In contrast to Ordinary Least Squares (OLS) that was used in the
250 original CCDC development, LASSO penalizes the sum of the absolute values of coefficients, in
251 some cases forcing a subset of the coefficients to zero. Together with the explicit limits enforced
252 on the number of coefficients, this reduces instances of overfitting, including in cases when
253 observations are too sparse or unevenly distributed in time to constrain the model to real
254 phenological features. To detect change, the LASSO model checks CCD model breaks with
255 respect to its last determined best-fit harmonic model.

256 To correctly detect change, the algorithm distinguishes between a substantive deviation from
257 model prediction and deviations that result from variability inherent in the data (due to
258 incomplete atmospheric removal and/or other sources of natural variation) to detect change. The
259 algorithm calculates two parameters related to dispersion, or scatter, to estimate the variability of
260 data for each spectral band. The first one is a comparison root-mean-square-error (RMSE) that is
261 the RMSE of the 24 observations covered by the model which are closest in day of year to the
262 last observation in the “peek window,” or over all observations covered by the model if there are
263 fewer than 24. This value is recalculated at each step of the time series. The second parameter
264 (*var*) is used to measure the overall variability of the data values and is defined as the median of
265 the absolute value of the differences between each observation and the *i*th successive
266 observation, where *i* is the smallest value such that the majority of these observation pairs are
267 separated by greater than 30 days, if possible (otherwise, *i*=1). The *var* is computed once at the
268 beginning of the standard procedure, using all non-masked observations in the time series.

269 Observations not yet incorporated into the model are evaluated as a group of no fewer than the
270 *PEEK_SIZE* parameter value; this is the “peek window,” which “slides” along the time series

271 one observation at a time. Each iteration, a value is calculated for each individual observation
272 within the peek window, as follows:

$$mag_n = \sum_{i \in D} \left(\frac{resid_{n,i}}{\max(var_i, RMSE_i)} \right)^2 \quad (2)$$

273 where, $resid_{n,i}$ is the residual relative to the LASSO models for each band i , for each
274 observation n within the *PEEK_SIZE* window, var_i and $RMSE_i$ are the parameters of dispersion
275 as described above, for each band i . This summation is carried out for all bands i in the set of
276 *DETECTION_BANDS* (D). This produces a scalar magnitude, representing the deviation from
277 model prediction across these bands, for each observation. The detection of a model break
278 requires this value to be above the *CHANGE_THRESHOLD* value for all observations in the
279 window. This is separate from the value that is reported as a per-band magnitude when a change
280 is detected in the time series. Change detection sensitivity depends on the value of change
281 threshold. The *CHANGE_THRESHOLD* is determined in Eqs. S2 and S3 in the Supplementary.
282 If $mag_n < CHANGE_THRESHOLD$ for any n in the *Peek_Size* window, then add the most
283 recent observation to the segment by shifting the *Peek_Size* window one observation forward in
284 the time series. If $mag_n > CHANGE_THRESHOLD$ for all n in the *Peek_Size* window, this is
285 considered a spectral break.

286 **3.3 Permanent snow and insufficient clear observation procedures**

287 The permanent snow procedure indicates that too few clear (less than 25% of total observations)
288 or water observations, which are identified from the QA band, exist to robustly detect change,
289 and a large fraction of observations are snow. The algorithm will return at most one segment that
290 fits through the entire time series and provide the filtered observations number at least twelve.
291 The model will, under the default settings, fit only four coefficients (i.e., characterizing the
292 reflectance and brightness temperature bands using only a simple harmonic with no higher
293 frequency terms). Unlike other procedures, snow pixels are not filtered out and are fit as part of
294 the annual pattern. This avoids overfitting the model to a seasonally sparse observation record.
295 Similarly, for the insufficient clear observations determined by the QA band, the model will
296 perform a LASSO regression fit for the entire time series using four coefficients. The model
297 coefficients and RMSE from this regression are recorded. Additional parameters including the
298 start, end, and observation count are also saved. Further, the change Boolean value is set to 0,

299 and the break day is recorded as the last observation date. The magnitude of change as zero for
300 each band is also saved.

301 **3.4 Land cover classification**

302 The CCDC algorithm characterizes the land cover component of a pixel at any point using the
303 LCMAP time series model approach from the Landsat 4–8 records. The classification of CCDC
304 is accomplished for every pixel based on data from the time series models (e.g., model
305 coefficients). Land cover classifications are generated on an annual basis, using July 1st as a
306 representative date. A list of land cover classes and descriptions is provided in Table 1. Fig.3
307 illustrates an overall classification approach.

308 **3.4.1 Classification algorithm**

309 We chose eXtreme Gradient Boosting (XGBoost) (Chen and Guestrin, 2016) as the classification
310 method. XGBoost is a scalable implementation of gradient tree boosting, which is a supervised
311 learning method that can be used to develop a classification model when provided with an
312 appropriate training dataset. Generally, for a given dataset, a tree ensemble model uses additive
313 functions, which correspond to independent tree structures, to predict the land cover. The
314 predictions from all trees are also normalized to the final class probabilities using the softmax
315 function. The algorithm can handle sparse data and theoretically justify weighted quantile sketch
316 for approximate learning. The resultant trained model can be applied to a larger dataset to
317 generate predictions and probability scores which are the basis for LCMAP primary and
318 secondary land cover types. The primary and secondary land cover confidence values are
319 calculated from these scores.

320 **3.4.2 Training dataset**

321 The training data used in XGBoost for the LCMAP Collection 1.0 land cover products is from
322 the USGS NLCD 2001 land cover product (Homer et al., 2020). To meet the LCMAP land cover
323 legend, the NLCD data is first cross-walked to LCMAP classes, as shown in Fig.4 and Table 2.
324 The use of NLCD data that was cross-walked to the LCMAP land cover legend as the training
325 data will reduce uncertainties and improve the consistency of annual land cover change. For
326 example, grass and shrub have different ecological functions. Their spectral signatures are
327 distinct in some ecological regions but are very close in others, especially in the western

328 ecoregions of the conterminous United States (Underwood et al., 2007; Xian et al., 2013). Grass
329 and shrub usually grow close together, making it difficult to separate them in thematic land
330 cover. Combining these two cover classes can reduce uncertainties potentially caused by lack of
331 spectral distinction in Landsat observations. Furthermore, the extent of each land cover in the
332 cross-walked NLCD layer is eroded by one pixel. This step aims to reduce potential noise in the
333 classifier by removing pixels that may be heavily mixed with different cover types, or whose
334 land cover label may be less reliable. It also removes the narrow linear low-intensity developed
335 pixels corresponding to road networks, which were found to have registration issues with
336 Landsat ARD in some areas.

337

338 **3.4.3 Ancillary data**

339 Ancillary data used in the classification contains two main datasets: the DEM and the WPI layer.
340 Three DEM derivative datasets are implemented as geographic references for land cover
341 classification as ancillary data including topographic slope, aspect, and position index. The WPI
342 is highly related to wetland distribution and has a potential to improve wetland classification in
343 LCMAP.

344 **3.4.4 Classification procedures**

345 For each pixel, CCD segment data for the segment that includes the July 1st, 2001 date is used
346 with training data to create classification models (Zhou et al., 2020; Zhu et al., 2016). Data
347 generated from the CCD models are used to make the land cover classification because different
348 land cover classes can have different shapes for the estimated time series models. The
349 coefficients of the CCD models including the overall mean and model coefficients except
350 intercepts can be used to estimate the intra-annual changes caused by phenology and sun angle
351 differences for the i th Landsat band. The information obtained from the time series model is
352 useful for land cover classification. The CCD model data used with training data include the
353 model coefficients (except the intercepts) generated from surface reflectance and brightness
354 temperature bands, the model RMSE value for each band, and an average intercept value that is
355 calculated from average annual reflectance values for each band for the July 1, 2001 year. The
356 model training procedure is conducted at the tile level, using random samples drawn from the
357 targeted tile as well as the eight surrounding tiles to avoid not having enough training samples of

358 rare land cover types in the targeted tile. Cross-walked and eroded NLCD data are used for
359 classification labels, while the CCD model outputs and ancillary data are provided as
360 independent variables. Based on training data testing using different sample sizes, a target
361 sample size of 20 million pixels from the extent of 3x3 ARD tiles is chosen, requiring
362 approximately proportional representation of classes with the added constraint that no class be
363 represented by fewer than 600,000 or more than 8 million samples. If there are fewer than
364 600,000 samples available for a class, then all of the available samples are used without any
365 oversampling. The XGBoost hyperparameters are selected as maximum tree depth 8, fast
366 histogram optimized approximate greedy algorithm for tree method, multiclass logloss for
367 evaluation metric, and maximum number of rounds 500.

368 After the classification models in a given tile are trained, predictions are generated for each July
369 1st date that has an associated CCD segment (Fig. 5). The prediction information is supplied to
370 the production step for the creation of land cover. The process is repeated for each tile for the
371 entire CONUS ARD extent.

372 **3.5 Validation data**

373 The LCMAP land cover product is validated using an independent reference dataset. The
374 reference data, which consists of 24,971 30 m x 30 m pixels selected via a simple random
375 sampling method over CONUS, is collected from these sample plots between 1985 and 2017.
376 The TimeSync tool is used to efficiently display Landsat data for interpretation and to record
377 these interpretations into a database (Cohen et al., 2010; Pengra et al., 2020b; Stehman et al.,
378 2021). TimeSync displays the input Landsat images in two basic ways: by annual time-series
379 images and by pixel values plotted through time. For the image display, single 255 x 255-pixel
380 subsets of Landsat images in the growing season are displayed in sequence from 1984 to 2018.
381 Trained interpreters have access to all available images in each year to collect attributes in three
382 basic categories: 1) land use, 2) land cover, and 3) change processes. Additional attribute details
383 for the change processes, such as clear-cut and thinning associated with harvest events, are also
384 collected. The interpreters manually label these attributes using Landsat 5, 7, and 8 imagery,
385 high-resolution aerial photography, and other ancillary datasets (Cohen et al., 2010; Pengra et al.,
386 2020b). Interpreters also use ancillary data to support interpretation of Landsat and high-
387 resolution imagery, although Landsat data takes the highest weight of evidence. Recording the

388 full set of attributes in land use, land cover, and land change categories provides sufficient
389 information to meet the needs of LCMAP as well as other potential users. Quality assurance and
390 quality control (QA/QC) processes are also implemented to ensure the quality and consistency of
391 the reference data among interpreters and over the time span of data collection (Pengra et al.,
392 2020b). Each reference sample is interpreted by a trained interpreter and about 60% of these
393 pixels are interpreted independently by a second analyst. Much of the QA/QC process relies on
394 comparing the interpretations at these duplicated sample pixels. Duplicated sample pixels that
395 have interpreter disagreement are evaluated in the QA/QC process, focusing on identifying
396 issues with specific classes or interpreters, flagging sample pixels for further review and possible
397 editing, and providing ongoing training and feedback to interpreters throughout the collection
398 process. QA/QC related reviews are also completed on sample pixels that show interpretation
399 data such as uncommon and/or illogical land use and land cover combinations, multi-year
400 disturbance processes, rare classes, or other opportunistically identified situations. Interpreted
401 attributes of sample pixels are edited, if necessary, to create the final attribute assignments for
402 the reference data. These final attributes are then cross-walked to a single LCMAP land cover
403 class label, providing a single land cover reference label for each year of the time series for each
404 sample pixel.

405 The validation analysis protocols focus on estimating the confusion matrix and overall, user's,
406 and producer's accuracy by comparing the reference data and product data labels. Overall
407 accuracy and producer's accuracy as well as standard errors are produced using post stratified
408 estimators (Card, 1982; Stehman, 2013). For accuracy estimates that are produced by combining
409 multiple years of data, the sampling design is treated as a one-stage cluster sample where each
410 pixel represents a cluster and each year of observation is the secondary sampling unit using
411 cluster sampling standard error formulas (Pengra et al., 2020b; Stehman et al., 2021). The
412 validation is only performed for primary land cover and change products, not for other LCMAP
413 science products (Supplementary Section 4).

414 **3.6 Information warehouse and data store**

415 LCMAP adopts an information warehouse and data store (IW+DS) system that can expand
416 storage solutions along with data access and discovery services running on the EROS Shared
417 Mesos Cluster. The system provides different storage solutions to allow for flexibility in

418 choosing what best fits a dataset’s characteristics and currently comprises Apache Cassandra
419 (<https://cassandra.apache.org/>) and Ceph (<https://ceph.io/>) object storage. The services provide
420 data ingest, retrieval, discovery, metadata, processing, and other functionalities. LCMAP
421 maintains a copy of Landsat Collection 1 ARD and other similarly tiled ancillary datasets that
422 are spatially subset within the IW+DS to allow efficient retrieval and to enable large-scale
423 CCDC processing and other algorithmic work. The ingest process is designed to avoid bringing
424 in ARD tile observations that are already present within the IW+DS, to keep the input consistent
425 with any prior usage while allowing CCDC to bring in new observations as they are available.
426 Algorithmic results, products, and other intermediate data are kept on the Ceph object store
427 arranged using a prefix structure to label the identity of the data, with the actual object names
428 incorporating spatial concepts such as tile and chip that is a small subset of a tile and contains
429 100 by 100 30-m pixels.

430

431 **4 Results and Discussion**

432

433 The LCMAP primary land cover and change products were evaluated to outline annual land
434 cover change from 1985 to 2017 in the conterminous Unites States.

435 **4.1 Collection 1.0 primary land cover distribution and change**

436 The CONUS primary land cover mapping result and the primary confidence in 2010 are shown
437 in Fig. 6a and b, respectively. The land cover map illustrates distributions of different land cover
438 types across CONUS. The primary confidence is above 90% for most land cover classes,
439 suggesting that the classification models were created with high confidence for land cover
440 mapping for most classes in most regions. Some vegetation transition (green in Fig. 6b) occurs
441 mainly in the southeast region suggesting gradual tree recovery from disturbances associated
442 with tree harvesting. Fig. 6c and d display numbers of land cover changes and spectral changes
443 detected by the CCDC model between 1985 and 2017. The number of land cover changes
444 represents how many times land cover has changed from one type to another for a specific pixel.
445 However, the number of spectral changes denotes how many times the model has detected
446 spectral changes in a CCD time series model where spectral observations have diverged from the
447 model predictions. These changes could relate to a change in thematic land cover or might

448 represent more subtle conditional surface changes. The southeast region shows more frequent
449 land cover changes in the 33 years (Fig. 6c). The western part of CONUS, however, contains
450 more spectral changes than in the east (Fig.6d). The NLCD land change estimates also show
451 similar change patterns between 2001 and 2016 (Homer et al., 2020). The different spatial
452 patterns in the total number of land cover changes (Fig. 6c) and detected spectral changes (Fig.
453 6d) suggest that not all changes lead to land cover change (e.g., drought and precipitation-related
454 changes in vegetation or grassland fire). The large numbers of spectral change were mainly
455 detected in the southern grassland area.

456 Fig. 7 shows the temporal changes of areas for eight land cover classes from 1985 to 2017.
457 Among all classes, grass/shrub, tree cover, and cropland were dominant land cover types,
458 followed by wetland, water, developed, barren, and snow/ice. The land cover and change
459 datasets show that developed land has a consistent increasing trend with an 8.4% increase while
460 barren increased 9.1% between 1985 and 2017. Overall, the developed and barren areas
461 increased $2.58 \times 10^4 \text{ km}^2$ and $8.56 \times 10^3 \text{ km}^2$, respectively. Other land cover categories do not have
462 such increasing patterns. As for water, although fluctuating, it had a generally increasing trend.
463 The area of wetland had a rapid decrease before 2000, following a relatively steady though
464 fluctuating trend. Net wetland extent declined about 0.4% from 1985 to 2017. The grass/shrub
465 and tree cover classes both experienced consistent increasing trends before 2008 and 1995,
466 respectively, with areas reaching about $2.85 \times 10^6 \text{ km}^2$ for grass/shrub and $2.14 \times 10^6 \text{ km}^2$ for tree
467 in these two years. These two land covers gradually decreased since then. Tree cover declines
468 after 1996, showing a decreasing rate of 2.8% between 1985 and 2017. The cropland decreased
469 from 1985 to 2008 and quickly increased after that. By 2017, the area of cropland reached a
470 similar level of cropland area in 1988. Furthermore, most land cover changes are located in the
471 southeast region where many pixels change more than one time. The changes detected by the
472 CCD model suggest that landscape in the Midwest and west are more dynamic than in the east.
473 Many areas experience multiple disturbances although most of these changes do not result in
474 land cover transition.

475 The south ARD tile outlined in Fig. 6(a) covers the northern Dallas region, and the spatial
476 patterns of land cover and change are shown in more detail in Fig. 8. The land cover distributions
477 in the region show that urban land expands considerably from 1985 (Fig. 8a), to 1990 (Fig. 8b),

478 and to 2016 (Fig. 8c). The land conversion was primarily from cropland and grass/shrub to
479 developed land. Lake Ray Roberts was created in the late 1980s and captured in the land cover
480 map (Fig. 8b&c). The lake and urban conversion are also visible in the change count from 1985
481 to 2016 (Fig. 8g), which mainly show as blue, suggesting a one-time conversion. On the other
482 hand, there is almost no change in the urban center (Fig. 8g). Fig. 8 (d-f) shows high
483 classification confidence at the urban center, water, grass/shrub, and tree cover areas, whereas
484 cropland has relatively low confidence, indicating frequent management activities over croplands
485 in the regions. The total pixels of different change numbers suggest that one to two change times
486 are dominant, although some pixels change more than three times (Fig. 8h). The land cover
487 distributions in 1985, 1990, and 2017 show an increase in developed land and decreases in
488 cropland and grass/shrub (Fig. 8i).

489 The spatial patterns of land cover and change in the north ARD tile displayed in Fig. 6(a) in
490 northern Wyoming are shown in Fig. 9. The tile covers most of Yellowstone National Park, in
491 which tree, grass/shrub, and water are three dominant land cover types. Land cover in 1985,
492 1990, and 2016 (Fig. 9a-c) changed from tree to grass/shrub and back to tree cover. The primary
493 land cover confidence layers exhibit changes as decreasing vegetation from tree to grass/shrub
494 and increasing vegetation from grass/shrub to tree (Fig. 9d-f). For those trees and water bodies
495 that did not experience any disturbances, their magnitudes of confidence are relatively large. The
496 change map suggests that most forest lands experienced at least one change and some areas
497 changed multiple times (Fig. 9g). Most changes in forest lands were related to wildland fires that
498 occurred in the region. In 1988, 50 fires burned a mosaic covering nearly 3213 km² in
499 Yellowstone as a result of extremely warm, dry, and windy weather (NPS, 2021). Trees regrew
500 in some of the burn areas and these changes could occur more than once as shown in the change
501 map that indicates at least two changes in these areas. The total pixels of different change
502 frequencies suggest that one to two changes were dominant and very few pixels changed more
503 than three times (Fig. 9h). The land cover distributions in 1985, 1990, and 2017 had increases in
504 grass/shrub after 1985 and reductions in tree cover after that (Fig. 9i).

505 **4.2 Validation of land cover product**

506 The overall accuracy between the annual reference land cover label and the LCMAP annual land
507 cover products was calculated as 82.5% ($\pm 0.22\%$, standard error) when summarized for all years.

508 Overall accuracy across the time series (1985-2017) varied within about 1.5% annually, ranging
509 from a high of 83% in the late 1990s to about 82% in the late 2010s (Fig. 10). Per class
510 accuracies across CONUS ranged between 43% and 96% for user's accuracy (Table 3), with
511 water showing the highest accuracy (96% \pm 0.5% user's accuracy and 93% \pm 0.7% producer's
512 accuracy). Cropland has about 93% (\pm 0.3%) producer's accuracy and 70% (\pm 0.6%) user's
513 accuracy. The lowest accuracies are observed for barren and wetland. The per class per year
514 agreements show the accuracies vary slightly for each class in each year (Table 4). The
515 variations of annual overall accuracy are within a range of about 1.5% across the time series. The
516 slight decline in annual overall accuracy suggests that year-to-year trends may be a result of a
517 complex interplay of temporal biases in the LCMAP algorithm, Landsat data quality and
518 quantity, the model break detection accuracy of the LCMAP CCD, and errors in the training data
519 used for the classification. For example, the change detection portion of the algorithm is known
520 to be conservative in identifying land cover change. The CCD model assumes that the spectral
521 variations of the land surface through time can be characterized with annual harmonic models
522 and can be separated into discrete periods of time. Therefore, the model performs better when the
523 short-term spectral variability of the land surface is low, the changes have a large spectral
524 response, and the observational data density is high. Over time, the actual land cover may evolve
525 away from the phenology represented by spectral models that may have missed one or more
526 spectral breaks, which will impact accuracy especially when the land cover changes are
527 persistent rather than cyclic, such as with an expanding urban footprint. Annual accuracy of
528 Developed showed an upward trend in user accuracy (UA) and a downward trend in producer
529 accuracy (PA) over time (Stehman et al., 2021). The increasing availability of high-resolution
530 data used by the interpreters may have increased the likelihood of identifying features
531 characteristic of Developed land that could not be identified earlier in the time series, leading to
532 an increase in the proportion of Developed area estimated from the sample. Consequently, the
533 increasing sensitivity of the reference interpretation to landscape features may account for the
534 difference between the mapping and the reference data over time. Lower data density toward the
535 beginning and end of the time series may decrease accuracy, which when combined with other
536 factors, can contribute to the annual land cover overall accuracy across all years.

537 **4.3 Significance of the product**

538 One of the biggest advances of LCMAP relative to conventional methods available to date is its
539 approach of generating annual land change products by using the entire Landsat archive at a
540 large geographic scale. Landsat ARD, which is the foundation for LCMAP, is effective and
541 straightforward for tracking and characterizing the historical land changes at a pixel level over
542 decades. Compared to conventional methods, detecting changes using all available observations
543 enables us to date these changes as they occur. After change is detected, temporally consistent
544 land cover products rather than stochastic changes in labels can be produced at annual intervals
545 by conducting classification from CCD model segmented contributions

546 The LCMAP product suite includes five land cover change and five land surface change science
547 products. It represents a new paradigm that consistently and continuously provides a large
548 volume of land change information for land change monitoring, land resource management, and
549 scientific research. In addition to primary and secondary land cover before and after changes,
550 change segments containing spectral change time and magnitude are provided to explore the
551 changes in land condition and could meet various user communities' needs. The LCMAP
552 products can improve our understanding of causes, rates, and consequences of the land surface
553 changes (Rover et al. 2020) such as forest changes caused by wildfire and insect outbreaks.

554 By implementing the CCDC algorithm through a system engineering approach, LCMAP
555 provides a fully automated framework for land change monitoring. The framework can also be
556 updated to include the latest Landsat records so that it can be used for operational continuous
557 monitoring in a large geographic extent (Brown et al. 2020). Therefore, when new observations
558 become available, the framework can provide timely and consistent land cover characteristics to
559 the public.

560 **4.4 Limitations and challenges**

561 Although LCMAP Collection 1.0 products have been proven to be successful in detecting
562 various land surface changes to support research applications related to environment and ecology
563 conditions, limitations and challenges exist. Utilizing Landsat ARD data as input provided
564 consistent time series Landsat imagery with high level geometric and radiometric quality for
565 implementing the CCDC method. Nevertheless, the densities of Landsat observation records
566 varied greatly across space and time due to spatial differences in Landsat scene overlap and
567 temporal coverage, as well as regional differences in contamination by clouds, cloud shadows,

568 and snow. The change detection accuracies of CCD models were highly influenced by the
569 temporal frequency of available observations. Zhou et al. (2019) found that using harmonized
570 Landsat-8 and Sentinel-2 (HLS) data increased the temporal frequency of the data and thus
571 enhanced the ability to model seasonal variation and derived better change detection results than
572 using Landsat data alone. Integrating multi-mission data could provide the opportunity to
573 enhance change detection, especially for the land cover types that are highly dynamic or in
574 frequently cloudy/snowy areas.

575 Providing only eight general land cover classes and their changes in LCMAP Collection 1.0
576 products limits the usage of the product in some applications that need a higher level of thematic
577 land cover detail. For example, shrub and grass are two major vegetation types and have
578 different ecological functions, but they are not delineated separately in LCMAP Collection 1.0
579 products. Lack of measurement of grassland-shrub transition constrains the study of shrub
580 encroachment, which is a symptom of land degradation. However, NLCD 2001 level I land
581 cover product had different mapping accuracies for different land cover types in different
582 ecological regions (Wickham et al., 2010). For example, the grass mapping accuracies were
583 higher in the eastern regions than they were in most western mapping regions. The accuracies of
584 shrub cover had similar variation patterns across CONUS. These accuracy variations suggest
585 uncertainties of the products, especially in most western regions where grass and shrub are more
586 difficult to be separated. Combining grass and shrub from the NLCD 2001 product reduced
587 uncertainties introduced by the two individual components and made the accuracy of the
588 grass/shrub product in LCMAP relatively high and consistent across CONUS (Stehman et al.,
589 2021). NLCD has established new efforts to improve mapping accuracies by adding innovative
590 approaches for land cover classification and introducing continuous rangeland products in
591 western CONUS for NLCD thematic land cover products since 2001 (Homer et al., 2020). The
592 use of new NLCD products as the training data will support LCMAP to produce more land cover
593 types including separating grass and shrub in the future.

594 Adopting NLCD land cover product as the training data source efficiently provided abundant
595 training samples to deliver land cover product with high classification accuracy. Selecting a
596 sufficient size of training samples is important for CCDC models to obtain accurate
597 classification. Previous land cover post-classification analysis suggested that the overall

598 classification accuracy increased when the training samples increased (Gong et al., 2020). The
599 recent global land cover classification also suggested that the appropriate training sample size for
600 a mapping extent of three 158 km x 158 km tiles should be larger than 60,000 (Zhang et al.,
601 2021). For the LCMAP land cover classification, a much larger training size was utilized to
602 ensure that these training samples could represent landscape features in the classification tiles.
603 However, these training data were randomly selected from the NLCD land cover product,
604 suggesting errors could potentially be carried over to the training samples due to potential errors
605 in the training source. Besides uncertainties in training data, some obvious challenges such as
606 class definitional differences between pasture/hay and grassland between NLCD and LCMAP
607 could potentially be carried over to the LCMAP land cover product. Improving training data by
608 reducing uncertainties and potential errors in a more consistent and accurate way is critical to
609 strengthen land cover classification and to improve the scientific quality of LCMAP products in
610 the future.

611 There are apparent shifts in some land cover types, especially in snow/ice and barren (Fig.7), and
612 a decline in overall agreement (Fig.10) in 2017, the last year of the Collection 1.0 product. The
613 last year's product usually is provisional because limited Landsat observations are available at
614 the end of a time series. The CCDC requires at least 24 clear observations to create full models
615 for change detection and classification. Without sufficient clear observations, the algorithm
616 could not produce model break accurately. Therefore, in the last year of a time series, the rule-
617 based assignment is implemented to label land cover for these pixels that do not have enough
618 observations to build a time series model. Both primary and secondary land cover classes are
619 assigned from the last identified primary and secondary classes.

620

621 **5 Data Availability**

622 The LCMAP products generated in this paper are available at <https://earthexplorer.usgs.gov/>
623 (LCMAP, 2021). All LCMAP land change products are mosaiced for the conterminous United
624 States in the GeoTIFF format. Find exact data as described here at
625 <https://doi.org/10.5066/P9W1TO6E>. The reference dataset used for the product validation is also
626 available at <https://www.sciencebase.gov/catalog/item/5e57e965e4b01d50924a93f6>
627 or <https://doi.org/10.5066/P98EC5XR> (Pengra et al., 2020a).

628

629 **6 Conclusions**

630 The continuous Landsat observations spanning from the 1980s to the present, new generations of
631 change detection and classification models, and systems capable of processing large volume data
632 are offering unprecedented opportunities to characterize land cover and detect land surface
633 change consistently and accurately. Additionally, the collection of reference data used to validate
634 land cover products provides validation result for each land cover category annually. To capture
635 the variability of landscape condition and its responses to different disturbances, land cover and
636 land surface change datasets need to be produced over a large geographic scale. LCMAP has
637 produced a suite of land change product at 30-m resolution including the reference dataset in the
638 United States. In that context, LCMAP was developed to generate an essential dataset to meet
639 broad scientific research and resource management needs. Using the CCDC algorithm and
640 Landsat ARD to determine whether change has occurred at any given point in the observation
641 record, LCMAP produced annual land cover and change datasets for the conterminous United
642 States in a robust manner. These new datasets and the novel production systems will allow for
643 new generation of research and applications in connecting time series remote sensing
644 observations with land surface change at a much finer scale than previously possible.

645

646 **Supplement.** The supplement related to this article is attached.

647

648 **Author contributions.**

649 KS conducted PyCCD programming for CCD/CCDC models. ZZ developed the original
650 MATLAB version of CCD/CCDC programs. JH participated in reference data collection. DW
651 and QZ assisted in data integration tasks. GX analysed the data and wrote the manuscript with
652 contributions from all co-authors.

653

654 **Completing interests.** The authors declare that they have no conflict of interest.

655

656 **Acknowledgements.**

657 Any use of trade, firm, or product names is for descriptive purposes only and does not imply
658 endorsement by the U.S. Government. Qiang Zhou and Congcong Li's work were performed
659 under Work performed under USGS contract 140G0119C0001.

660

661

662

References

- Anderson, J. R., Hardy, E.E., Roach, J.T., and Witmer, R.E.: A land use and land cover classification system for use with remote sensor data, Geological Survey Professional Paper, 964, 1-28, 1976.
- Boryan, C., Yang, Z., Mueller, R., and & Craig, M.: Monitoring US agriculture: the US department of agriculture, national agricultural statistics service, cropland data layer program, Geocarto International, 26, 341-358, 2011.
- 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., Pengra, B. W., Smith, K., Horton, J. A., Xian, G., Auch, R. F., Sohl, T. L., Sayler, K. L., Gallant, A. L., Zelenak, D., Reker, R. R., and Rover, J.: Lessons learned implementing an operational continuous United States national land change monitoring capability: The Land Change Monitoring, Assessment, and Projection (LCMAP) approach, Remote Sensing of Environment, 238, 111356, 2020.
- Bullock, E. L., Woodcock, C. E., and Holden, C. E.: Improved change monitoring using an ensemble of time series algorithms, Remote Sensing of Environment, 238, 2020.
- Chen, J., Liao, A., Cao, X., Chen, L., Chen, Z., He, C., Han, G., Peng, S., Lu, M., and Zhang, W.: Global land cover mapping at 30 m resolution: A POK-based operational approach, ISPRS journal of photogrammetry and remote sensing : official publication of the International Society for Photogrammetry and Remote Sensing, 103, 7-27, 2015.
- Chen, T. and Guestrin, C.: XGBoost, Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 785-794, 2016.
- Cohen, W. B., Yang, Z., and & Kennedy, R.: Detecting trends in forest disturbance and recovery using yearly Landsat time series: 2. TimeSync — Tools for calibration and validation, Remote Sensing of Environment, 114, 2911-2924, 2010.
- Dwyer, J. L., Roy, D. P., Sauer, B., Jenkerson, C. B., Zhang, H. K., and Lymburner, L.: Analysis Ready Data: Enabling Analysis of the Landsat Archive, Remote Sensing, 10, 1363, 2018.
- Erb, K. H., Luysaert, S., Meyfroidt, P., Pongratz, J., Don, A., Kloster, S., Kuemmerle, T., Fetzel, T., Fuchs, R., Herold, M., Haberl, H., Jones, C. D., Marin-Spiotta, E., McCallum, I., Robertson, E., Seufert, V., Fritz, S., Valade, A., Wiltshire, A., and Dolman, A. J.: Land management: data availability and process understanding for global change studies, Glob Chang Biol, 23, 512-533, 2017.
- Foley, J. A., DeFries, R., Asner, G.P., Barford, C., Bonan, G., Carpenter, S.R., Chapin, F.S., Coe, M.T., Daily, G.C., Gibbs, H.K., Helkowski, J.H., Holloway, T., Howard, E.A., Kucharik, C.J., Monfreda, C., Patz, J.A., Colin Prentice, I., Ramankutty, N., Synder, P.K.: Global consequences of land use, Science, 309, 570-574, 2005.
- Foley, J. A., Ramankutty, N., Brauman, K. A., Cassidy, E. S., Gerber, J. S., Johnston, M., Mueller, N. D., O'Connell, C., Ray, D. K., West, P. C., Balzer, C., Bennett, E. M., Carpenter, S. R., Hill, J., Monfreda, C., Polasky, S., Rockstrom, J., Sheehan, J., Siebert, S., Tilman, D., and Zaks, D. P.: Solutions for a cultivated planet, Nature, 478, 337-342, 2011.
- Franklin, S. E., Ahmed, O. S., Wulder, M. A., White, J. C., Hermosilla, T., and Coops, N. C.: Large Area Mapping of Annual Land Cover Dynamics Using Multitemporal Change Detection and Classification of Landsat Time Series Data, Canadian Journal of Remote Sensing, 41, 293-314, 2015.
- Friedl, M. A., Sulla-Menashe, D., Tan, B., Schneider, A., Ramankutty, N., Sibley, A., and Huang, X.: MODIS Collection 5 Global Land Cover: Algorithm Refinements and Characterization of New Datasets, Remote Sensing of Environment, 114, 168-182, 2010.

Gong, P., Li, X., Wang, J., Bai, Y., Chen, B., Hu, T., Liu, X., Xu, B., Yang, J., Zhang, W., and Zhou, Y.: Annual maps of global artificial impervious areas (GAIA) between 1985 and 2018, *Remote Sensing of Environment*, 236, 111510, 2020.

Gong, P., Liu, H., Zhang, M., Li, C., Wang, J., Huang, H., Clinton, N., Ji, L., Li, W., Bai, Y., Chen, B., Xu, B., Zhu, Z., Yuan, C., Ping Suen, H., Guo, J., Xu, N., Li, W., Zhao, Y., Yang, J., Yu, C., Wang, X., Fu, H., Yu, L., Dronova, I., Hui, F., Cheng, X., Shi, X., Xiao, F., Liu, Q., and Song, L.: Stable classification with limited sample: transferring a 30-m resolution sample set collected in 2015 to mapping 10-m resolution global land cover in 2017, *Science Bulletin*, 64, 370-373, 2019.

Hansen, M. C., Potapov, P.V., Moore, R., et al.: High-resolution global maps of 21st century forest cover change, *Science*, 342, 850-853, 2013.

Hermosilla, T., Wulder, M. A., White, J. C., Coops, N. C., and Hobart, G. W.: Disturbance-Informed Annual Land Cover Classification Maps of Canada's Forested Ecosystems for a 29-Year Landsat Time Series, *Canadian Journal of Remote Sensing*, 44, 67-87, 2018.

Homer, C., Dewitz, J., Jin, S., Xian, G., Costello, C., Danielson, P., Gass, L., Funk, M., Wickham, J., Stehman, S., Auch, R., and Riitters, K.: Conterminous United States land cover change patterns 2001–2016 from the 2016 National Land Cover Database, *ISPRS Journal of Photogrammetry and Remote Sensing*, 162, 184-199, 2020.

Jin, S., Yang, L., Danielson, P., Homer, C., Fry, J., and Xian, G.: A comprehensive change detection method for updating the National Land Cover Database to circa 2011, *Remote Sensing of Environment*, 132, 159-175, 2013.

Kennedy, R. E., Yang, Z., Braaten, J., Copass, C., Antonova, N., Jordan, C., and Nelson, P.: Attribution of disturbance change agent from Landsat time-series in support of habitat monitoring in the Puget Sound region, USA, *Remote Sensing of Environment*, 166, 271-285, 2015.

LCMAP, LCMAP Collection 1 Science Products [data]. <https://doi.org/10.5066/P9W1TO6E>, 2021.

Li, X., Zhou, Y., Meng, L., Asrar, G. R., Lu, C., and Wu, Q.: A dataset of 30 m annual vegetation phenology indicators (1985–2015) in urban areas of the conterminous United States, *Earth System Science Data*, 11, 881-894, 2019.

Li, X., Zhou, Y., Meng, L., Asrar, G. R., Lu, C., and Wu, Q.: A dataset of 30 m annual vegetation phenology indicators (1985–2015) in urban areas of the conterminous United States, *Earth System Science Data*, 11, 881-894, 2019.

Li, X., Zhou, Y., Zhu, Z., and Cao, W.: A national dataset of 30 m annual urban extent dynamics (1985–2015) in the conterminous United States, *Earth System Science Data*, 12, 357-371, 2020.

NPS, 2021. Fire - Yellowstone National Park, <https://www.nps.gov/yell/learn/nature/fire.htm#:~:text=Number%20in%20Yellowstone,human%20caused%20fires%20were%20suppressed.&text=The%20number%20of%20fires%20has,70%20C285%20acres%20in%20Yellowstone%20burned>. Accessed in April 27, 2021.

Pengra, B. W., Stehman, S. V., Horton, J. A., and Wellington, D. F.: Land Change Monitoring, Assessment, and Projection (LCMAP) Version 1.0 Annual Land Cover and Land Cover Change Validation Tables., U.S. Geological Survey data release, [data], doi: <https://doi.org/10.5066/P98EC5XR>, 2020a.

Pengra, B. W., Stehman, S. V., Horton, J. A., Dockter, D. J., Schroeder, T. A., Yang, Z., and Loveland, T. R.: Quality control and assessment of interpreter consistency of annual land cover

reference data in an operational national monitoring program, *Remote Sensing of Environment*, 238, 111261, 2020b.

Picotte, J. J., Dockter, D., Long, J., Tolk, B., Davidson, A., and Peterson, B.: LANDFIRE remap prototype mapping effort: Developing a new framework for mapping vegetation classification, change, and structure, *Fire*, 2, 35, 2019.

Reid, W. V., Chen, D., Goldfarb, L., Hackmann, H., Lee, Y. T., Mokhele, K., Ostrom, E., Raivio, K., Rockstrom, J., Schellnhuber, H. J., and Whyte, A.: Earth System Science for Global Sustainability: Grand Challenges, *Science*, 330, 916-917, 2010.

Stehman, S. V., Pengra, B. W., Horton, J. A., and Wellington, D. F.: Validation of the U.S. Geological Survey's Land Change Monitoring, Assessment and Projection (LCMAP) Collection 1.0 annual land cover products 1985–2017, *Remote Sensing of Environment*, 265, 2021.

Szantoi, Z., Geller, G. N., Tsendbazar, N.-E., See, L., Griffiths, P., Fritz, S., Gong, P., Herold, M., Mora, B., and Obregón, A.: Addressing the need for improved land cover map products for policy support, *Environmental Science & Policy*, 112, 28-35, 2020.

Tibshirani, R.: Regression shrinkage and selection via the lasso, *Journal of the Royal Statistical Society: Series B (Methodological)*, 58, 267-288, 1996.

Turner II, B. L., Lambin, E. F., and Reenberg, A.: The emergence of land change science for global environmental change and sustainability, *Proceedings of the National Academy of Sciences of the United States of America*, 104, 20666-20671, 2007.

Underwood, E. C., Ustin, S. L., and Ramirez, C. M.: A comparison of spatial and spectral image resolution for mapping invasive plants in coastal California, *Environ Manage*, 39, 63-83, 2007.

Wickham, J. D., Stehman, S. V., Fry, J. A., Smith, J. H., and Homer, C. G.: Thematic accuracy of the NLCD 2001 land cover for the conterminous United States, *Remote Sensing of Environment*, 114, 1286-1296, 2010.

Wulder, M. A., Coops, N. C., Roy, D. P., White, J. C., and Hermosilla, T.: Land cover 2.0, *International Journal of Remote Sensing*, 39, 4254-4284, 2018.

Xian, G., Homer, C., Meyer, D., and Granneman, B.: An approach for characterizing the distribution of shrubland ecosystem components as continuous fields as part of NLCD, *ISPRS Journal of Photogrammetry and Remote Sensing*, 86, 136-149, 2013.

Zhou, Q., Tollerud, H. J., Barber, C. P., Smith, K., and Zelenak, D.: Training data selection for annual land cover classification for the land change monitoring, assessment, and projection (LCMAP) initiative, *Remote Sensing*, 12, 699, 2020.

Zhu, Z., Gallant, A. L., Woodcock, C. E., Pengra, B., Olofsson, P., Loveland, T. R., Jin, S., Dahal, D., Yang, L., and R.F., A. A.: Optimizing selection of training and auxiliary data for operational land cover classification for the LCMAP initiative, *ISPRS Journal of Photogrammetry and Remote Sensing* 122, 206-221, 2016.

Zhu, Z., Wang, S., and Woodcock, C. E.: Improvement and expansion of the Fmask algorithm: Cloud, cloud shadow, and snow detection for Landsats 4–7, 8, and Sentinel 2 images, *Remote Sensing of Environment*, 159, 269-277, 2015a.

Zhu, Z. and Woodcock, C. E.: Automated cloud, cloud shadow, and snow detection in multitemporal Landsat data: An algorithm designed specifically for monitoring land cover change, *Remote Sensing of Environment*, 152, 217-234, 2014a.

Zhu, Z. and Woodcock, C. E.: Continuous change detection and classification of land cover using all available Landsat data, *Remote Sensing of Environment*, 144, 152-171, 2014b.

Zhu, Z. and Woodcock, C. E.: Object-based cloud and cloud shadow detection in Landsat imagery, *Remote Sensing of Environment*, 118, 83-94, 2012.

Zhu, Z., Woodcock, C. E., Holden, C., and Yang, Z.: Generating synthetic Landsat images based on all available Landsat data: Predicting Landsat surface reflectance at any given time. , Remote Sensing of Environment, 162, 67-83, 2015b.

Caption of Table

Table 1 LCMAP land cover product specifications

Table 2 NLCD land cover cross-walked to LCMAP land cover

Table 3. Confusion matrix for CONUS (all years combined) where cell entries represent percent of CONUS area. Overall accuracy is 82.5% ($\pm 0.22\%$). Standard errors for user's and producer's accuracies are shown in parentheses and n is the number of sample pixels for each row and column.

Table 4 Overall per class agreement in percentage between 1985 and 2017

Caption of Figure

Figure 1 Landsat ARD tile grids for the conterminous U.S.

Figure 2 Overall procedures of the CCD algorithm.

Figure 3 The overall approach of land cover classification in CCDC.

Figure 4. NLCD 2001 land cover (a), cross-walked LCMAP land cover classes (b), LCMAP land cover eroded by one pixel (c), zoomed in cross-walked land cover from NLCD 2001 (d), and zoomed in LCMAP land cover classes eroded by one pixel (e). The color legends represent NLCD land cover class and LCMAP primary land cover (LCPRI).

Figure 5 CCD change detection and segmentation using Landsat blue, green, red, near-infrared, short-wave infrared (SWIR) 1, short-wave infrared (SWIR) 2, and thermal bands. Blue dots are all available clear Landsat records in each year. The horizontal lines in different colors represent land cover classes labeled by the algorithm. The vertical lines show model break dates. The back line is the model fits. The high-resolution images show landscape conditions in 2007 and 2013.

Figure 6 Illustration of the LCMAP product: (a) Primary land cover in 2010, (b) Primary land cover confidence in 2010, (c) the frequency of land cover changes from 1985 to 2017, and (d) total number of spectral changes detected from 1985 to 2017.

Figure 7 Areal variations of eight primary land cover types from 1985 to 2017 in CONUS.

Figure 8 Primary land cover and confidences in 1985 (a) and (d), 1990 (b) and (e), 2016(c) and (f), change in 1985-2017 (g), the frequency of land cover change (x-axis) from 1985 to 2017 and numbers of total pixels (y-axis) of these changes of different change (h), and areas (y-axis) of different land cover (x-axis) in the three times for the ARD tile 16_14 (i).

Figure 9 Primary land cover and confidences in 1985 (a) and (d), 1990 (b) and (e), 2016 (c) and (f), and change in 1985-2017 (g), the frequency of land cover change (x-axis) from 1985 to 2017 and numbers of pixels (y-axis) of these changes (h), and areas (y-axis) of different land cover (x-axis) in the three times for the ARD tile 9_6 (i).

Figure 10 Overall agreement between LCMAP primary land cover and reference data across CONUS. The cross lines represent \pm one standard errors.

Table 1 LCMAP land cover product specifications

Code	Land Cover Class	Description
1	Developed	Areas of intensive use with much of the land covered with structures (e.g., high-density residential, commercial, industrial, mining, or transportation), or less intensive uses where the land cover matrix includes vegetation, bare ground, and structures (e.g., low-density residential, recreational facilities, cemeteries, transportation/utility corridors, etc.), including any land functionality related to the developed or built-up activity.
2	Cropland	Land in either a vegetated or unvegetated state used in production of food, fiber, and fuels. This includes cultivated and uncultivated croplands, hay lands, orchards, vineyards, and confined livestock operations. Forest plantations are considered as forests or woodlands (Tree Cover class) regardless of the use of the wood products.
3	Grass/Shrub	Land predominantly covered with shrubs and perennial or annual natural and domesticated grasses (e.g. pasture), forbs, or other forms of herbaceous vegetation. The grass and shrub cover must comprise at least 10% of the area and tree cover is less than 10% of the area.
4	Tree Cover	Tree-covered land where the tree cover density is greater than 10%. Cleared or harvested trees (i.e. clearcuts) will be mapped according to current cover (e.g. Barren, Grass/Shrub).
5	Water Bodies	Areas covered with water, such as streams, canals, lakes, reservoirs, bays, or oceans.
6	Wetland	Lands where water saturation is the determining factor in soil characteristics, vegetation types, and animal communities. Wetlands are composed of mosaics of water, bare soil, and herbaceous or wooded vegetated cover.
7	Ice and Snow	Land where accumulated snow and ice does not completely melt during the summer period (i.e. perennial ice/snow).
8	Barren	Land comprised of natural occurrences of soils, sand, or rocks where less than 10% of the area is vegetated.

Table 2 NLCD land cover cross-walked to LCMAP land cover

NLCD Value	LCMAP Value
Water	Water
Ice/Snow	Ice and Snow
Developed, open space; Developed, low intensity; Developed medium intensity; Developed, high intensity	Developed
Barren	Barren
Deciduous forest, Evergreen forest, Mixed forest	Tree Cover
Shrub/Scrub, Grassland/Herbaceous	Grass/Shrub
Hay/Pasture, Cultivated crops	Cropland
Woody wetland, Emergent herbaceous wetland	Wetland

Table 3. Confusion matrix for CONUS (all years combined) where cell entries represent percent of CONUS area. Overall accuracy is 82.5% ($\pm 0.22\%$). Standard errors for user's and producer's accuracies are shown in parentheses and n is the number of sample pixels for each row and column.

Map	Devel	Crop.	Grass /Shrub	Tree	Water	Wetland	Ice/ Snow	Barren	Total	User (SE)	n
Devel.	3.000	0.139	0.321	0.377	0.024	0.035		0.001	3.896	77 (1.2)	32102
Crop.	0.918	16.527	5.061	0.799	0.027	0.368		0.003	23.702	70 (0.6)	195283
Grass /Shrub	0.368	0.757	30.649	2.599	0.045	0.229		0.332	34.980	88 (0.3)	288197
Tree	0.340	0.143	1.414	23.387	0.049	0.579		0.006	25.917	90 (0.3)	213531
Water	0.013	0.008	0.048	0.024	4.788	0.067		0.020	4.968	96 (0.5)	40932
Wetland	0.062	0.129	0.361	0.944	0.172	3.688		0.001	5.357	69 (1.3)	44136
Ice/Sno w			0.004	0.004		0.004	0.012	0.004	0.028	43 (18.7)	231
Barren	0.072	0.005	0.501	0.013	0.056	0.012		0.492	1.151	43 (2.8)	9485
Total	4.772	17.707	38.358	28.149	5.162	4.981	0.012	0.859	100.00		
Prod (SE)	63 (1.3)	93 (0.3)	80 (0.4)	83 (0.4)	93 (0.7)	74 (1.2)	100 (0)	57 (3.2)			
n	39319	145886	316027	231916	42530	41042	99	7078			

Table 4 Overall per class agreement in percentage between 1985 and 2017

Overall Per Class Agreement	Developed	Cropland	Grass/Shrub	Tree	Water	Wetland	Snow/Ice	Barren
1985	66	80	83	87	95	72	60	49
1986	67	80	83	87	95	72	60	49
1987	68	80	83	86	95	72	60	49
1988	68	80	83	87	95	72	60	49
1989	68	80	84	87	95	72	60	48
1990	68	80	84	87	95	72	60	48
1991	68	80	84	87	95	72	60	49
1992	69	80	84	87	95	71	60	50
1993	69	80	84	87	95	71	60	49
1994	69	80	84	87	95	71	60	49
1995	70	80	84	87	95	72	60	49
1996	69	80	84	87	95	72	60	48
1997	70	80	84	87	95	72	60	49
1998	70	80	84	87	94	72	60	48
1999	70	80	84	87	95	72	60	48
2000	70	80	84	87	95	72	60	48
2001	70	80	84	87	95	72	60	49
2002	70	80	84	86	95	72	60	49
2003	70	80	84	87	94	71	60	48
2004	69	80	84	86	94	71	60	48
2005	70	80	84	86	94	71	60	49
2006	70	79	84	86	94	71	60	49
2007	70	79	84	86	94	71	60	50
2008	70	79	84	86	94	71	60	49
2009	70	79	84	86	94	71	60	49
2010	70	79	84	86	94	71	60	50
2011	70	79	84	86	94	71	60	51
2012	70	79	83	86	94	71	60	50
2013	69	79	83	86	94	71	60	50
2014	69	79	83	86	94	71	60	50
2015	69	79	83	86	94	71	60	50
2016	69	79	83	86	94	71	60	50
2017	69	78	83	85	94	70	60	49

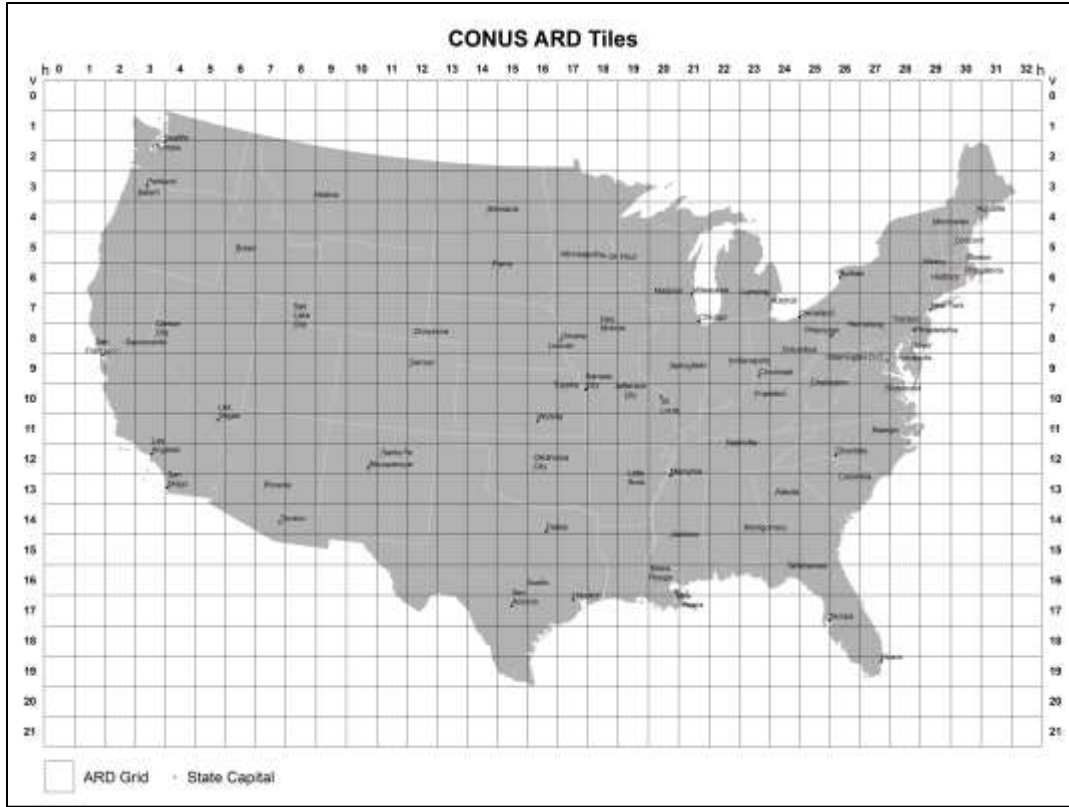


Figure 1 Landsat ARD tile grids for the conterminous U.S.

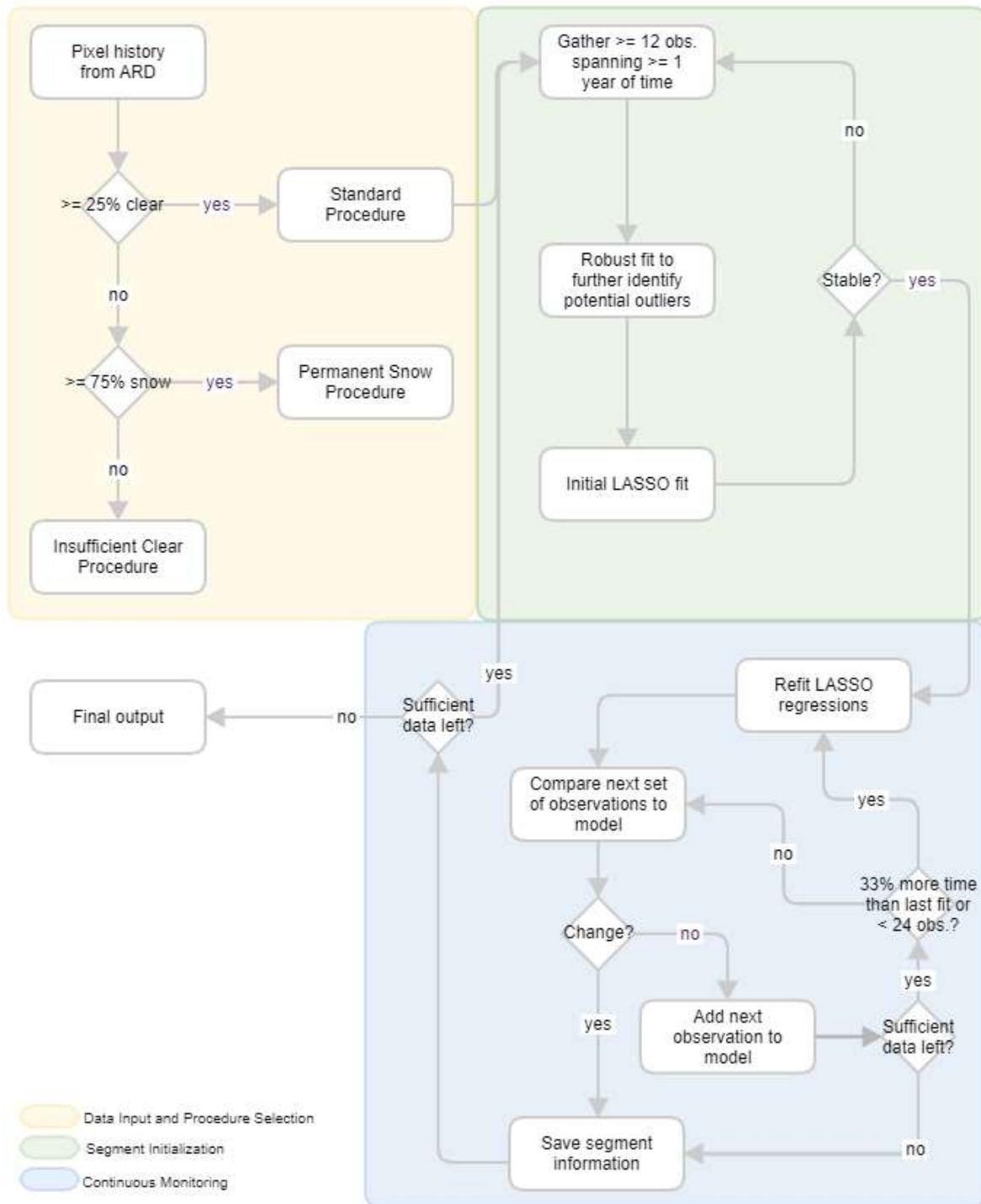


Figure 2 Overall procedures of the CCD algorithm.

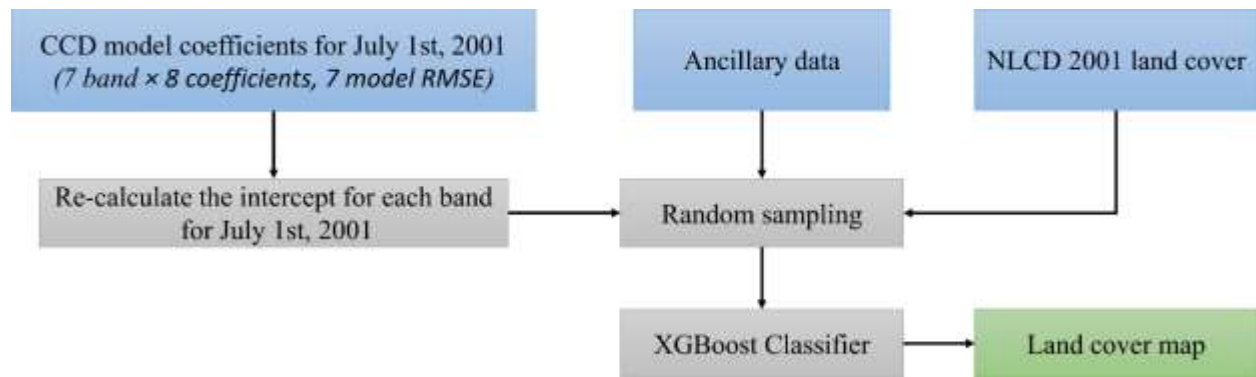


Figure 3 The overall approach of land cover classification in CCDC.

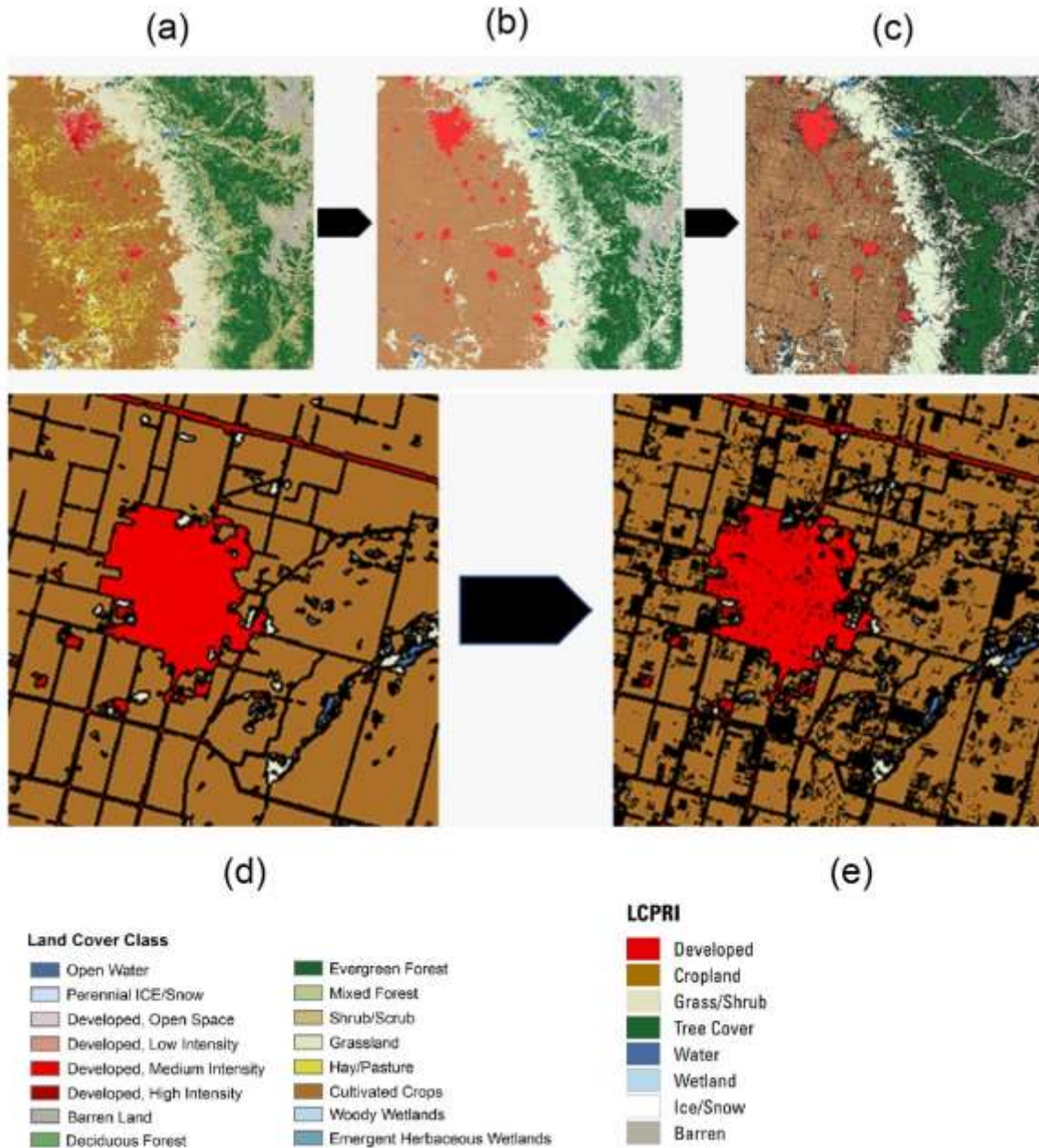


Figure 4. NLCD 2001 land cover (a), cross-walked LCMAP land cover classes (b), LCMAP land cover eroded by one pixel (c), zoomed in cross-walked land cover from NLCD 2001 (d), and zoomed in LCMAP land cover classes eroded by one pixel (e). The color legends represent NLCD land cover class and LCMAP primary land cover (LCPRI).

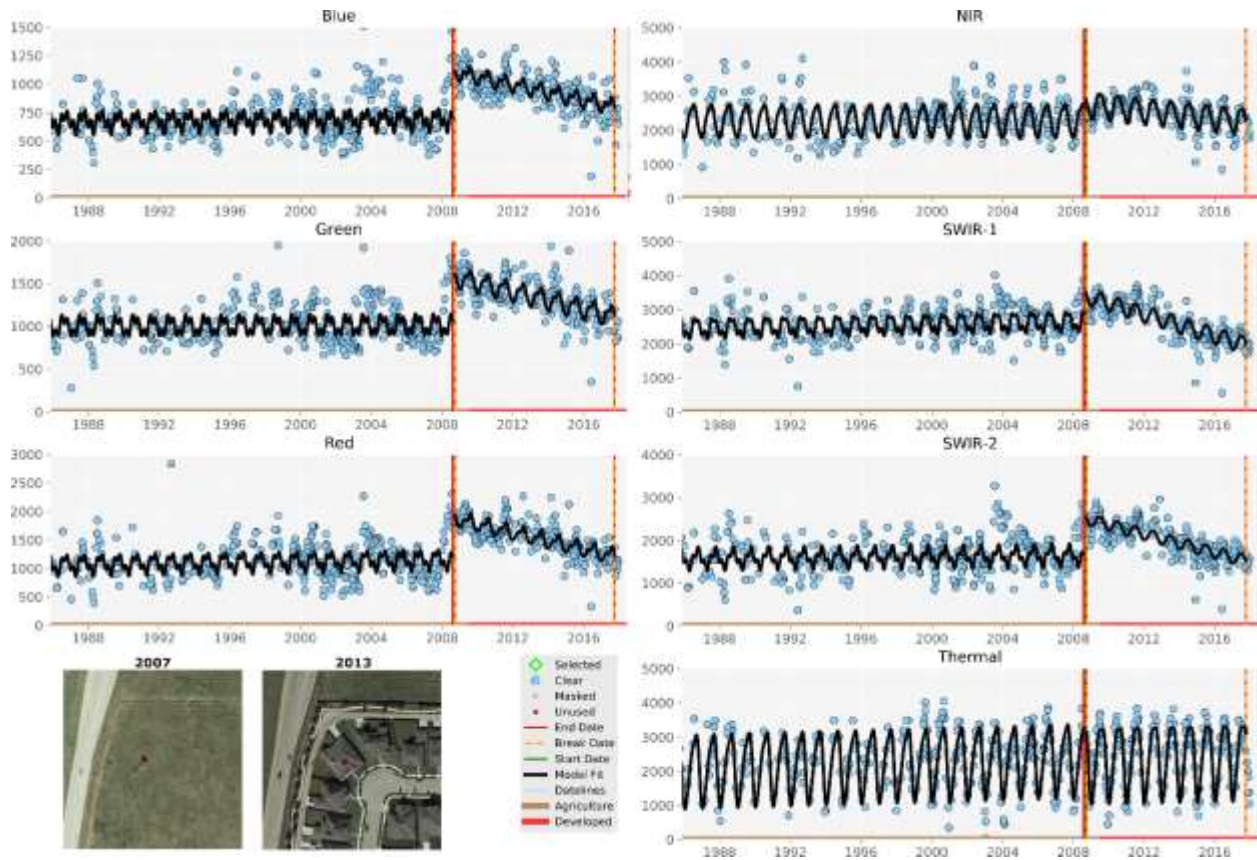


Figure 5 CCD change detection and segmentation using Landsat blue, green, red, near-infrared, short-wave infrared (SWIR) 1, short-wave infrared (SWIR) 2, and thermal bands. Blue dots are all available clear Landsat records in each year. The horizontal lines in different colors represent land cover classes labeled by the algorithm. The vertical lines show model break dates. The back line is the model fits. The high-resolution images show landscape conditions in 2007 and 2013.

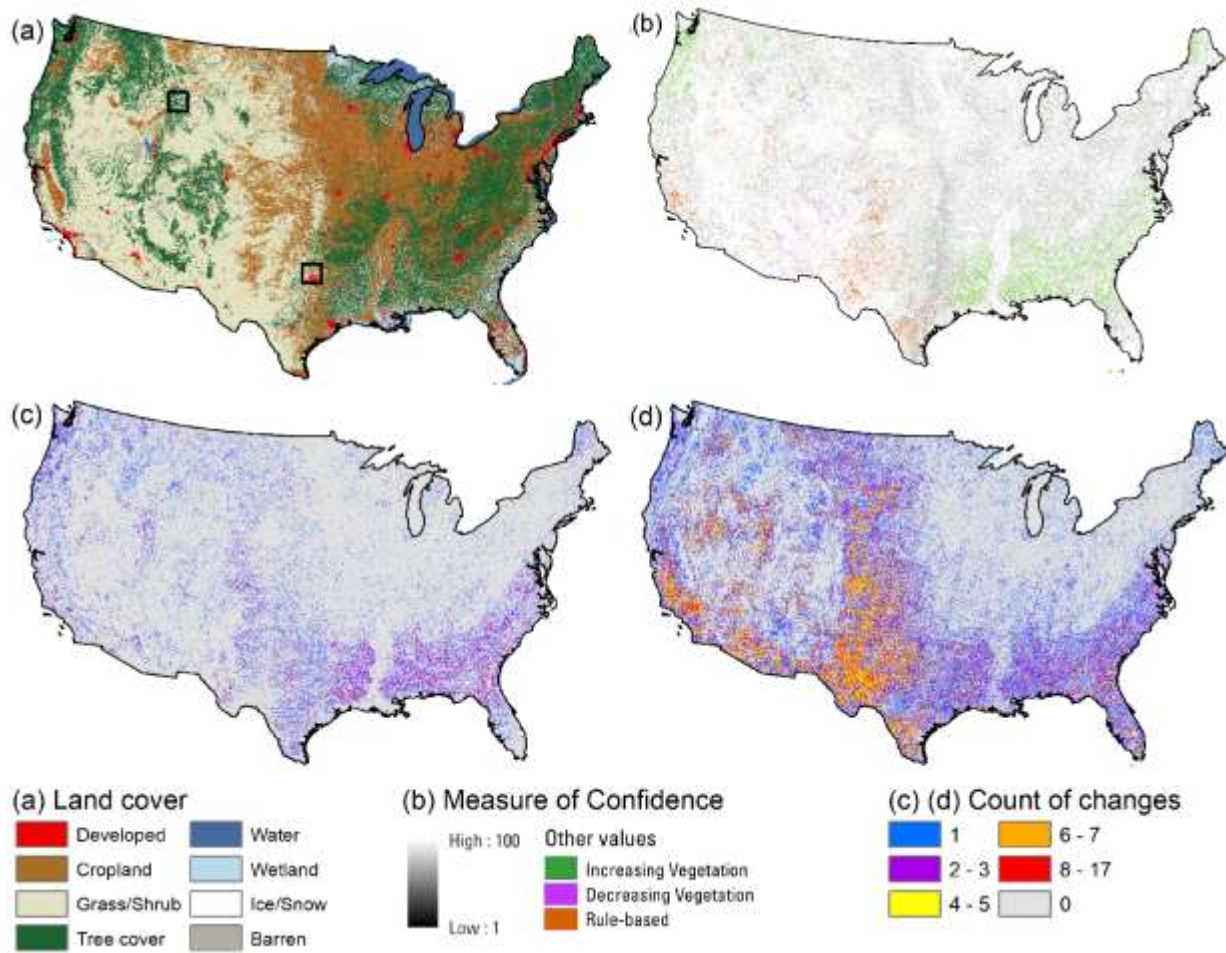


Figure 6 Illustration of the LCMAP product: (a) Primary land cover in 2010, (b) Primary land cover confidence in 2010, (c) the frequency of land cover changes from 1985 to 2017, and (d) total number of spectral changes detected from 1985 to 2017.

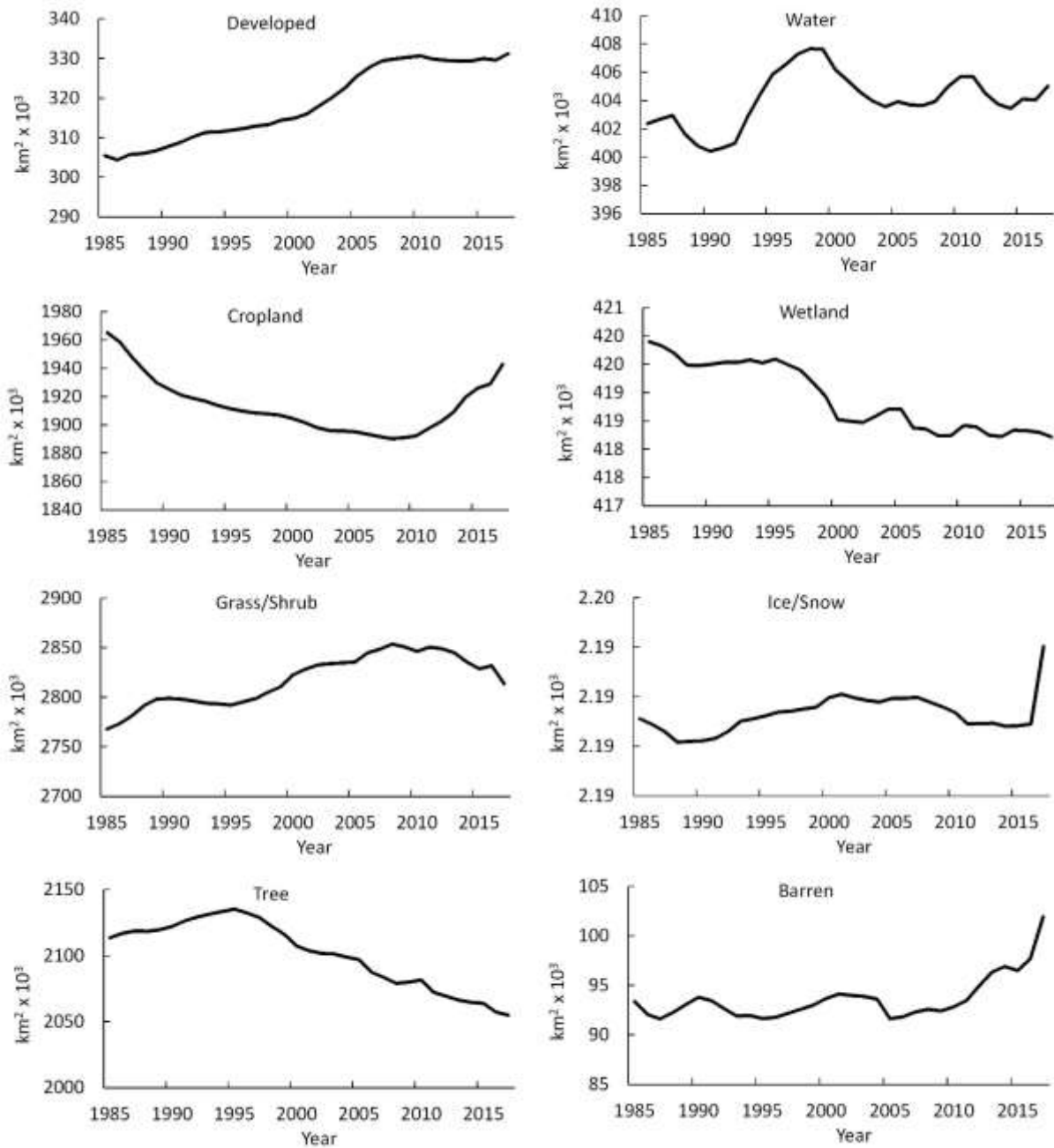


Figure 7 Areal variations of eight primary land cover types from 1985 to 2017 in CONUS.

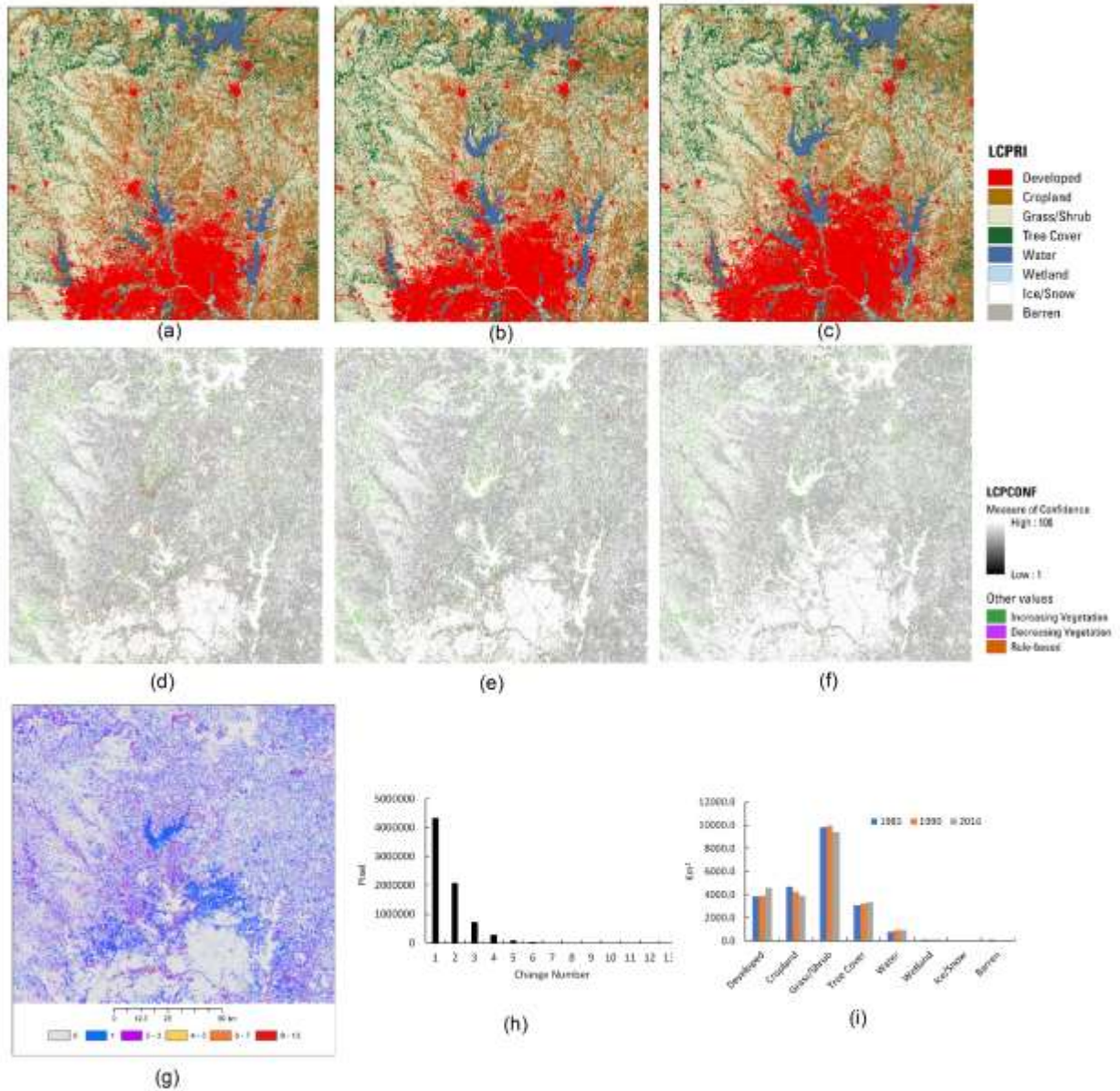


Figure 8 Primary land cover and confidences in 1985 (a) and (d), 1990 (b) and (e), 2016(c) and (f), change in 1985-2017 (g), the frequency of land cover change (x-axis) from 1985 to 2017 and numbers of pixels (y-axis) of these changes (h), and areas (y-axis) of different land cover (x-axis) in the three times for the ARD tile 16_14 (i).

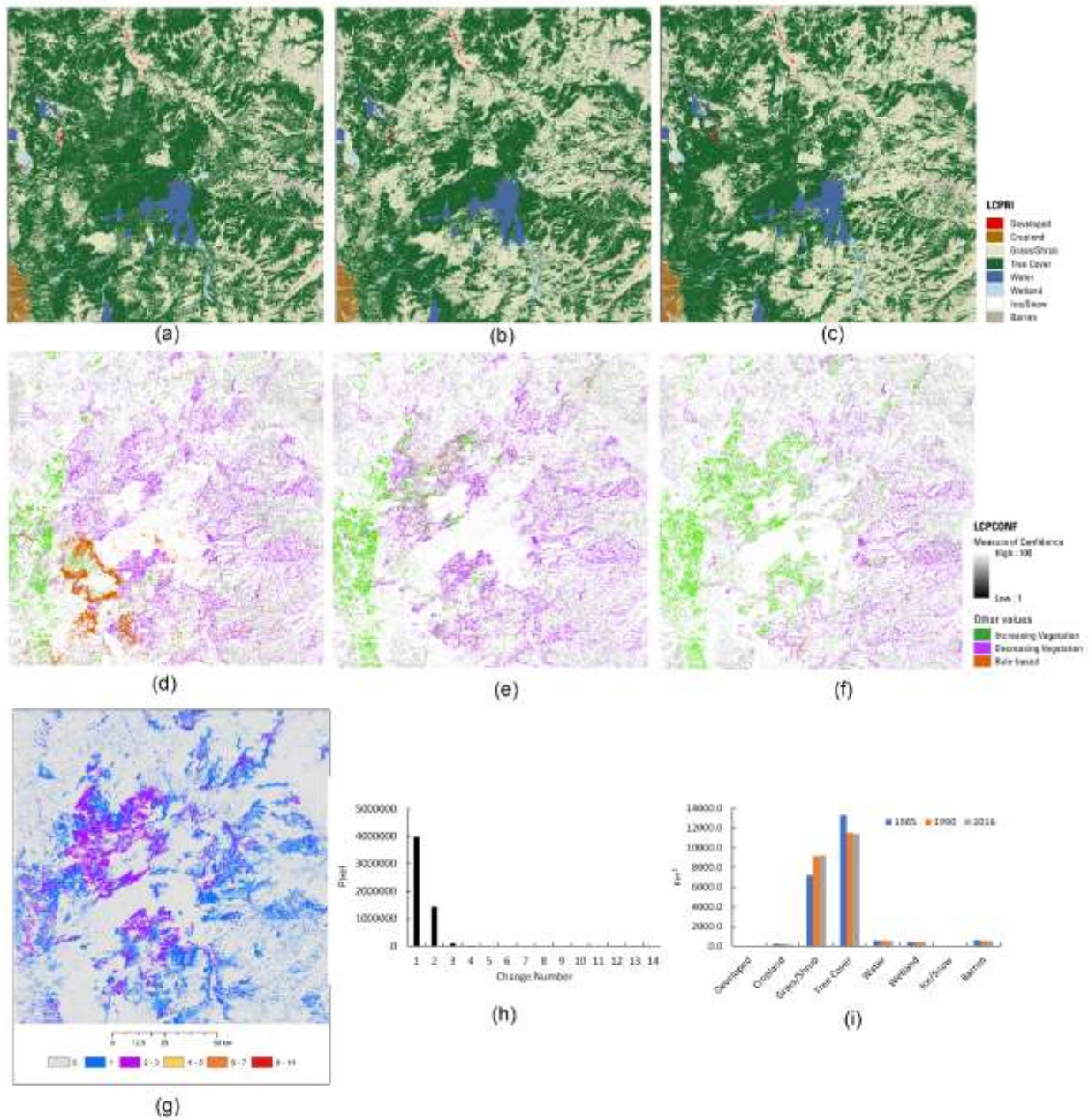


Figure 9 Primary land cover and confidences in 1985 (a) and (d), 1990 (b) and (e), 2016 (c) and (f), and change in 1985-2017 (g), the frequency of land cover change (x-axis) from 1985 to 2017 and numbers of pixels (y-axis) of these changes (h), and areas (y-axis) of different land cover (x-axis) in the three times for the ARD tile 9_6 (i).

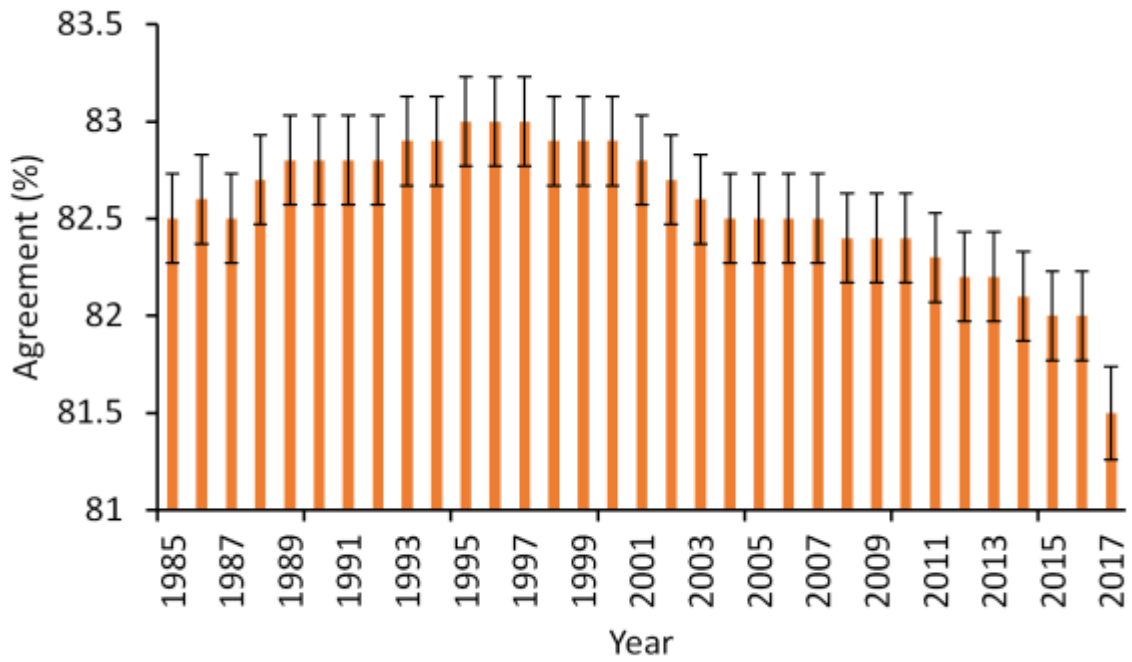


Figure 10 Overall agreement between LCMAP primary land cover and reference data across CONUS. The cross lines represent +/- one standard errors.