

## Continuous and consistent land use/cover change estimates using socio-ecological data

Michael Marshall<sup>a†</sup>, Michael Norton-Griffiths<sup>a</sup>, Harvey Herr<sup>a</sup>, Richard Lamprey<sup>b</sup>, Justin Sheffield<sup>c</sup>, Tor Vagen<sup>a</sup>, and Joseph Okotto-Okotto<sup>d</sup>

<sup>a</sup> Climate Research Unit, World Agroforestry Centre, United Nations Ave, Gigiri, P.O. Box 30677-00100, Nairobi, Kenya, Email: m.marshall@cgiar.org

<sup>b</sup> Fauna & Flora International, The David Attenborough Building, Pembroke St, Cambridge, CB2 3QZ, UK

<sup>c</sup> Department of Civil and Environmental Engineering, Princeton University, Princeton, NJ, 08544, USA

<sup>d</sup> Lake Basin Development Authority, P.O Box 1516-40100, Kisumu, Kenya

<sup>†</sup> Corresponding Author

1

### 2 **Abstract**

3       A growing body of research shows the importance of land use/cover change (LULCC) on  
4 modifying the earth system. Land surface models are used to stimulate land-atmosphere  
5 dynamics at the macro- (regional to global) scale, but model bias and uncertainty remain that  
6 needs to be addressed, before the importance of LULCC is fully realized. In this study, we  
7 propose a method of improving LULCC estimates for land surface modelling exercises. The  
8 method is driven by projectable socio-ecological geospatial predictors available seamlessly  
9 across sub-Saharan Africa and yielded continuous (annual) estimates of LULCC at 5x5 km<sup>2</sup>  
10 spatial resolution. The method was developed with 2,252 5x5 km<sup>2</sup> sample area frames of the  
11 proportion of several land cover types in Kenya over multiple years. Forty-three socio-  
12 ecological predictors were evaluated for model development. Machine learning was used for  
13 data reduction and simple (functional) relationships defined by generalized additive models were  
14 constructed on a subset of the highest ranked predictors ( $p \leq 10$ ) to estimate LULCC. The  
15 predictors explained 62% and 65% of the variance in the proportion of agriculture and natural  
16 vegetation, respectively, but were less successful at estimating more descriptive land cover types.

1 In each case, population density on an annual basis was the highest ranked predictor. The  
2 approach was compared to a commonly used remote sensing classification procedure, given the  
3 wide use of such techniques for macro-scale LULCC detection, and out-performed it for each  
4 land cover type. The approach was used to demonstrate significant trends in expanding  
5 (declining) agricultural (natural vegetation) land cover in Kenya from 1983-2012, with the  
6 largest increases (declines) occurring in densely populated high agricultural production zones.  
7 Future work should address the improvement (development) of existing (new) geospatial  
8 predictors and issues of model scalability and transferability.

9 **Key words: earth system; land use/cover change; population-environment; remote sensing;**  
10 **vegetation function**

1 **1. Introduction**

2 Land use/cover change (LULCC) is an important concern for global environmental  
3 sustainability, because it can adversely affect surface albedo and heating (Davin and de Noblet-  
4 Ducoudré, 2010); evapotranspiration and other components of the hydrologic cycle (Sterling et  
5 al., 2013); local to regional climate with the coupling or indirect recycling of surface moisture  
6 (Makarieva et al., 2013); global climate via carbon and other greenhouse gas emissions  
7 (Anderson-Teixeira and DeLucia, 2011; Ward et al., 2014); and ecosystem services worsened by  
8 these impacts (Turner et al., 2013). Land surface models, which can be coupled to a regional or  
9 global climate model, are used to simulate land-atmosphere interactions retrospectively or  
10 prospectively (Pitman, 2003) to identify intervention “hotspots” or develop realistic land  
11 management scenarios at the macro- (regional to global) scale (Turner et al., 2007).  
12 Traditionally, spatially-explicit LULCC was not an input to land surface models, but was instead  
13 represented by structural (e.g. leaf area index) or physiological (e.g. stomatal resistance) changes  
14 in vegetation. LULCC was then mapped in parallel to characterize these changes. These early  
15 attempts have been replaced by fully coupled LULCC and land surface models (e.g. Shevliakova  
16 et al., 2009; Lawrence et al., 2012). Although the impact of LULCC on the earth system is well  
17 established and quantifiable, studies remain few, due in part to the inadequacy of LULCC  
18 estimates (Pielke et al., 2011). In order to further land-atmosphere interaction research, LULCC  
19 models must be developed that provide consistent estimates over long historical time frames,  
20 regular (annual) intervals, and large spatial domains at 5x5 km<sup>2</sup> spatial resolution; are projectable  
21 50-100 years into the future; and use a consistent classification approach (Meiyappan et al.,  
22 2014; Rounsevell et al., 2014; Verburg et al., 2011).

1 Heistermann et al. (2006) reviews the two primary categories of macro-scale LULCC  
2 models (geographic and economic) while Schaldach and Priess (2008) and Rounsevell et al.  
3 (2014) include reviews of blended or integrated approaches. The Conversion of Land Use and  
4 its Effects (CLUE) model (Veldkamp and Fresco, 1996; Verburg et al., 2002) is an example of a  
5 geographic technique. It identifies important socio- (population, economy, society, politics and  
6 planning, culture, and technology) ecological (climate, vegetation, soil, topography, and  
7 hydrology) predictors from observed LULCC data, which are related to each other statistically,  
8 and then cellular automata are used to simulate competition between the predicted LULC types  
9 and neighboring grid cells based on these relationships. Decision rules are typically used  
10 iteratively to guarantee realistic LULC transitions occur. LandSHIFT (Alcamo et al., 2011) is an  
11 example of an economic approach, because supply (LULC) is distributed on a grid cell basis by  
12 demand. Supply is determined from national estimates of crop yield and the net primary  
13 productivity of grasslands. Multi-criteria analysis, which involves applying cost functions and  
14 LULC constraints based on socio-ecological inputs, is used to define demand hierarchically and  
15 disaggregate supply over baseline or projected periods. Integrated approaches (e.g. CLUMondo:  
16 van Asselen and Verburg, 2013) are becoming more common, because they more adequately  
17 account for LULCC processes and the interaction of demand and trade with supply than  
18 economic or geographic models, respectively. Like most geographic and economic models,  
19 however, integrated models have a sound theoretical basis, but can be difficult to employ on a  
20 grid-cell basis at high spatial resolution at the macro-scale, because of data inconsistencies and  
21 incongruities and model complexity that can propagate error, as well as, the time and other  
22 resources needed to operate them. Earth observation (remote sensing) models are an important  
23 sub-category of the geographic approach, because they overcome many of these challenges,

1 making their operational use on grid-cell basis at high spatial resolution at the macro-scale more  
2 feasible.

3 Hansen and Loveland (2012) and Ban et al., (2015) present recent reviews of macro-scale  
4 remote sensing-based LULCC modeling. Remote sensing approaches use multivariate statistical  
5 techniques to classify land cover types based on the spectral or textural characteristics of gridded  
6 satellite data (DeFries et al., 1995). These approaches are simpler than integrated approaches,  
7 because they tend to capture change at a single resolution directly with no interaction between  
8 adjacent pixels. Remote sensing approaches, therefore, tend to be more parsimonious than  
9 integrated approaches and require less time for processing. Early remote sensing approaches  
10 involved daily coarse spatial resolution (8km) Advanced Very High Resolution Radiometer  
11 (AVHRR) data available from 1981. Large disagreement and uncertainties in the models, due to  
12 mixed pixel effects from small LULC patch size, as well diverse classification systems and  
13 methods, limited their use at the macro-scale in the past (Lepers et al., 2005). Recently,  
14 improved computational storage and processing and consensus on classification has facilitated  
15 the creation of consistent global LULCC maps at Landsat (30 m) resolution (Giri et al., 2013).  
16 GlobeLand30 (Chen et al., 2015), for example, uses a pixel-object-knowledge-based approach to  
17 classify Landsat images from spectrally-derived vegetation indices globally in 2000 and 2010.  
18 The use of Landsat data alone poses serious challenges to modeling LULCC on an annual basis:  
19 persistent cloud cover and a 16-day revisit cycle, makes retrieval of cloud-free pixels difficult;  
20 the Landsat platforms have been retired (Landsat 5), have failed (Landsat 6), suffer from  
21 technical problems (Landsat 7), or are only recently active (Landsat 8). To improve the temporal  
22 resolution and continuity of classification, other remote sensing products, such as the Global  
23 Forest Change product (Hansen et al., 2010), fuse Moderate-resolution Imaging

1 Spectroradiometer (MODIS) data available every 1-2 days at 250-500m spatial resolution with  
2 Landsat data. But these products are only available over the MODIS era (2000-present), making  
3 long-term classification difficult. In short, the major drawback of remote sensing approaches is  
4 that the temporal range and continuity necessary for long-term annual global change detection  
5 are often sacrificed for high ( $\leq 500\text{m}$ ) spatial resolution. Finally, remote sensing data is not  
6 projectable like other socio-ecological data, such as population density, precipitation, or  
7 temperature, limiting their use to retrospective analyses.

8         The purpose of this study was to propose a simple (functional) way to map LULCC at the  
9 macro-scale at  $5 \times 5 \text{ km}^2$  spatial resolution on an annual basis using socio-ecological predictors  
10 that are available on an annual basis and projectable 50-100 years into the future to facilitate  
11 land-atmosphere modeling and research. The method was developed using sample area frames  
12 consisting of continuous land cover proportions developed from multi-year aerial and ground  
13 surveys in Kenya over a 30-year period. The approach was compared with remote sensing  
14 predictors that have been used to classify land cover types based on their phenology. Kenya is  
15 an ideal location to develop such a method, because like many countries in sub-Saharan Africa  
16 (SSA) data is scarce compared to the Global North, and the impact of land modification on  
17 people and the environment is high (Lambin et al., 2003). In addition: 1) population density is  
18 highest in the most agriculturally productive areas due to unequitable land distribution and poor  
19 infrastructure (Jayne and Muyanga, 2012) making ecological determinants that are generally  
20 used to map LULCC potentially less relevant (Pricope et al., 2013); 2) agriculture is the primary  
21 source of livelihood and crops are mostly rainfed (Ngetich et al., 2014); and 3) inter-annual  
22 rainfall variability is high and frequently causes devastating droughts and floods (Held and  
23 Soden, 2006).

## 1 **2. Data and Methods**

### 2 *2.1 Study area*

3 Aerial surveys were conducted in 1983, 1985, 2012, and 2013, to assess changes in land  
4 cover over parts of the Lake Victoria basin and central region of Kenya (Machakos and Makueni  
5 areas). The surveys yielded 2,252 5x5 km<sup>2</sup> sample area frames covering 28,150 km<sup>2</sup> or  
6 approximately 47% of Kenya's arable lands (**Figure 1**). Olofsson et al. (2012) has suggested  
7 that 5x5 km<sup>2</sup> sample area frames are appropriate for evaluating macro-scale LULCC models.  
8 The lakeshore and lowlands of Lake Victoria basin are primarily tropical with one long rain  
9 season that extends from February to September (UNEP, 2008). The neighboring highlands  
10 follow a bimodal pattern and annual totals are higher than near the lakeshore, due to warm moist  
11 westerlies during the West African monsoon and orographic uplift. Central Kenya is drier and  
12 has two distinct rain seasons: long rains (March-June) and short monsoon rains (October-  
13 December). The Machakos area, which includes Muranga', Kiambu, and the northern part of  
14 Machakos, is humid subtropical and therefore wetter than Makueni to the south-east, which is  
15 semi-arid.

16 The probability (proportion) of various land cover types within each frame were available  
17 at two levels of specificity: level one (agriculture, natural vegetation, urban, and miscellaneous)  
18 and level two (crops, fallow, shrubs, savanna, wetlands, forest, and agroforestry). These two  
19 levels of specificity were analyzed to determine the level of detail that can be captured using  
20 coarse resolution geospatial data. Continuous data were used, because at 5x5 km<sup>2</sup> resolution,  
21 spatial heterogeneity makes discrete classification impractical. Agriculture included  
22 agroforestry, defined here as trees on a farm; crops (banana, coffee, maize, sugar cane, tea,  
23 wheat, and others); and pasture/fallow. Natural vegetation included savanna, shrubs (open and

1 closed), wetlands (perennial and permanent), and forest (evergreen and deciduous). Urban  
2 included built up structures, such as roads, homes, and towns. Miscellaneous included fish  
3 ponds and other water bodies, exposed rock, and charcoal pits. The frames were developed  
4 using an aerial point sampling approach (Norton-Griffiths, 1988): several thousand geotagged  
5 aerial photos were taken over parallel transects spaced 1 km apart at approximately 488 m  
6 (height-above-ground) in 1983/1985 and then again in 2012/2013, resulting in approximately 7  
7 aerial natural color analogue photos per frame with a ground-sampling-distance of < 1 cm in  
8 1983/1985 and 5 aerial natural color digital photos per frame with a ground-sampling-distance of  
9 6.5 cm in 2012/2013. The retrieval dates are shown in **Table 1**. A team of six technicians  
10 interpreted the photos on a rolling basis to minimize potential bias and errors that can occur from  
11 manual classification by different interpreters and for different years. The proportion of each  
12 land cover type (0-100%) was determined by manually classifying a grid of 320 randomly  
13 distributed points superimposed over each photo. For each year, all land cover types were  
14 represented and classified, but not all frames were interpreted and classified (**Figure 1**). The  
15 interpretations were validated via site visits and meetings with community stakeholders. The  
16 estimates were then averaged over the photos across interpreters to get the proportions for each  
17 frame. Further details on the 1983/1985 and 2012/2013 campaigns can be found in EcoSystems  
18 Ltd (1983), EcoSystems Ltd (1987), and Lamprey (2013).

## 19 *2.2 Macro-scale data handling and processing*

20 The development of the functional relationships from the sample area frames involved  
21 four major steps illustrated in **Figure 2**. Non-remote sensing and remote sensing predictors were  
22 selected after an exhaustive online search that are freely and seamlessly available across SSA, so  
23 that the relationships can be used in future studies across the continent for retrospective or



1 prospective analyses. Given the large number of predictors collected, machine learning was used  
2 to identify a subset of the most powerful predictors before constructing the functional  
3 relationships. The functional relationships were then evaluated against remote sensing predictors  
4 with hold-out samples and finally, used to demonstrate how the relationships can be used to  
5 reconstruct LULCC estimates continuously through time.

6 Forty-three non-remote sensing (climatic, hydrologic, socioeconomic, and topographic)  
7 and sixteen remote sensing (phenological) predictors of land cover change were compared and  
8 subset for model-building with the sample area frames. Either slowly-changing (long-term  
9 average/one-time value) or dynamic predictors were considered. The slowly-changing predictors  
10 and their sources are shown in **Table 2**. Using these predictors alone could streamline the  
11 modeling process. However, in reality, phenology, climate, and population change frequently, so  
12 these predictors were derived on an annual basis as well. The handling and processing of  
13 annually-changing or dynamic predictors are discussed in Sections 2.2.1-2.2.3. For the  
14 remainder of the paper, dynamic predictors include a “.d” extension. All of the geospatial data  
15 was projected to Africa Equidistant Conic (m) to facilitate distance calculations. The predictors  
16 were resampled to the finest resolution data (90x90 m<sup>2</sup>) and aggregated to 5x5 km<sup>2</sup> resolution for  
17 model-building.

## 18 2.2.1 Climate

19 Bioclimatic (BIOCLIM: Hijmans et al., 2005) variables were used to capture climatic  
20 differences in land cover types because they 1) provide biologically meaningful information and  
21 2) have been projected mid-21<sup>st</sup> century at high spatial resolution for SSA (AFRICLIM: Platts et  
22 al., 2014). Two additional climate parameters were included in the analysis, because they are  
23 potentially relevant and part of the Platts et al. (2014) dataset: atmospheric demand for moisture

1 (Potential Evapotranspiration- PET) and the Moisture Index. The BIOCLIM variables were  
2 computed on an annual basis from 1983-2012 using monthly temperature, shortwave incoming  
3 radiation, and precipitation. The variables were estimated using the “biovars” function in the  
4 “dismo” package in R (Hijmans et al., 2015). As with the Platts et al. (2014) dataset, PET was  
5 estimated using Hargreaves and Samani (1985).

6         The temperature/radiation and precipitation predictors were taken from the Princeton  
7 University high resolution meteorological forcing (PHF) (Chaney et al., 2014) and the Climate  
8 Hazards Group InfraRed Precipitation with Stations (CHIRPS) (Funk et al., 2014) datasets,  
9 respectively. PHF originally spanned 1979-2008, but was extended to 2012 for this study. It is a  
10 downscaled version of the Princeton University global meteorological forcing (PGF) dataset  
11 (Sheffield et al., 2006) for SSA. It assimilates new observation data, specifically station data  
12 from the U.S. National Climatic Data Center (NCDC) Integrated Surface Database (ISD) and has  
13 undergone more rigorous correction than the global dataset. PHF is a blend of the most up-to-  
14 date observation-based, remote sensing, and reanalysis data sources: the National Centers for  
15 Environmental Prediction–National Center for Atmospheric Research (NCEP-NCAR) reanalysis,  
16 Global Precipitation Climatology Project, Tropical Rainfall Measuring Mission (TRMM), the  
17 Climatic Research Unit (CRU), and the Surface Radiation Budget. The data is downscaled using  
18 elevation. The dataset includes precipitation, minimum/maximum temperature, pressure,  
19 shortwave and longwave radiation, specific humidity, and wind speed at a daily time step and  
20  $0.1^\circ$  (~10 km at the equator) resolution. CHIRPS is available at pentad (5-day) intervals and  
21  $0.05^\circ$  (~ 5km at the equator) spatial resolution from 1981-2012. Like PHF, CHIRPS is a blend  
22 of several observation-based, remote sensing, and reanalysis sources: geostationary thermal  
23 infrared satellite observations from the Climate Prediction Center and National Climatic Data

1 Center; TRMM; and NOAA-NCAR. CHIRPS was selected as the precipitation data source over  
 2 PHF, because it incorporates the largest collection of ground-based precipitation data in East  
 3 Africa and bias-correction is performed using the Climate Hazards Precipitation Climatology  
 4 (Funk et al., 2015).

### 5 2.2.2 Population density

6 Population density was derived from the UNEP/GRID-Sioux Falls African Population  
 7 Distribution Database (APDD) on an annual basis from 1983-2012. APDD consists of  
 8 population density at a spatial resolution of 2.5 arc-minutes  $^{\circ}$  ( $\sim 5 \times 5 \text{ km}^2$  at the equator) for base  
 9 years 1960, 1970, 1980, 1990, and 2000. The grids are derived from population statistics at  
 10 various administrative (district, province, etc.) levels and temporal scales, depending on the  
 11 availability of national population statistics. A detailed description of the derivation of gridded  
 12 population can be found in Deichmann (1996). Each grid cell represents “population potential”,  
 13 based on its proximity to the transportation network (roads, railroads, and navigable rivers, and  
 14 major towns/cities). Population at a given administrative level is then disaggregated according to  
 15 the population potential. Grid cells that are closer to the network have higher coefficients and  
 16 therefore receive a larger proportion of the population than grid cells further away. The base  
 17 years are then extrapolated with an exponential growth/decay function (Davis, 1995). For  
 18 consistency, the same function was used to distribute population between base years on an  
 19 annual basis for each grid cell:

$$\mathbf{P}_{i,j,t} = \mathbf{P}_{i,j,T} e^{\Delta t \mathbf{k}_{i,j}} \quad (2)$$

$$\mathbf{k}_{i,j} = \ln(\mathbf{P}_{T+10n} / \mathbf{P}_{T+10(n-1)}) / 10 \quad (3)$$

20  $\mathbf{P}_{i,j,t}$  is the interpolated population/population density for a given year ( $t$ ) and at grid cell  $i, j$ ,  $\mathbf{P}_{i,j,T}$   
 21 is the population/population density for a given base year (period = 10 years),  $\Delta t$  is the change in

1 time from the base year to the year being interpolated, and  $k_{i,j}$  (**Equation 2**) is the growth/decay  
 2 coefficient. The growth/decay coefficient is defined by  $P_{T+10(n-1)}$  (initial base year for iteration  $n$ )  
 3 and  $P_{T+10n}$  (last base year for iteration  $n$ ). The denominator was set to ten, because  $k_{i,j}$  accounted  
 4 for decadal trends. After 2000, population statistics were extrapolated to 2012 using the 1990-  
 5 2000 growth/decay coefficients.

### 6 2.2.3 Remote Sensing Predictors

7 The National Aeronautics and Space Administration's Global Inventory Modeling and  
 8 Mapping Studies (GIMMS) Normalized Difference Vegetation Index (NDVI) Version 3  
 9 (NDVI3g) (Pinzon and Tucker, 2014) was used to estimate the remote sensing predictors. NDVI  
 10 is a ratio-based vegetation index derived from Earth observation (AVHRR) surface reflectance in  
 11 the visible red and near infrared (NIR). NDVI approaching one (zero) is indicative of dense  
 12 vegetation (bare soil). NDVI3g is available at  $0.08^\circ$  ( $\sim 8 \times 8 \text{ km}^2$  at the equator) spatial resolution  
 13 and at a 15-day timestep from 1983-2013. NDVI3g has been compared to other long-term global  
 14 vegetation records and is considered the most appropriate for trend analyses (Tian et al., 2015).

15 The predictors were derived from NDVI using harmonic regression (Eastman et al.,  
 16 2009) on an annual basis from 1983-2012. Linear harmonic regression estimates the amplitude  
 17 (maximum) and phase (timing) of a fitted time series, but unless higher order harmonics are  
 18 introduced, linear harmonic regression is too rigid to account for outliers and multimodal  
 19 regimes commonly found in the tropics. To overcome these obstacles, non-linear harmonic  
 20 regression (Carrão et al., 2010) was used to estimate five phenological predictors:

$$\mathbf{NDVI}_{i,j,T} = \mathbf{M}_{i,j} + \mathbf{A}_{i,j} \cos(\omega_0 t + \phi + \alpha \cos(\omega_0 t + \varphi)) \quad (1)$$

21 Where  $\mathbf{NDVI}_{i,j,T}$  is NDVI3g at grid cell  $i, j$  and over period  $T$ , which in this case was 24,  
 22 because non-linear harmonic regression was computed on an annual basis from the 15-day data;

1 **M** is the intercept (annual mean NDVI); **A** is the amplitude;  $\phi$  is the annual phase; and  $\alpha$  and  $\varphi$   
2 are non-linear terms defining the strength of non-linearity (asymmetry) and non-linear phase  
3 (deceleration/acceleration of asymmetry), respectively. The frequency ( $\omega_0$ ) equals  $2\pi/T$ . The  
4 approach can be reduced to a linear harmonic oscillator by setting  $\alpha\cos(\omega_0t+\phi)$  to zero. The  
5 non-linear predictors were derived at each grid cell using the “nlsLM” function in the  
6 “minipack.lm” package in R (Elzhov et al., 2015). nlsLM uses the Levenberg-Marquardt  
7 optimization method (Moré, 1978) to find the non-linear least-squares fit. The function was  
8 constrained by the seed and boundary conditions described in Carrão et al., (2010). One  
9 thousand iterations at each grid cell were performed to avoid fitting local optima. Linear terms  
10 (**A** and  $\phi$ ) were computed for the analysis as well, using the “lm” function in the “stats” package  
11 in R (<https://cran.r-project.org/>), because they are more efficient and are easier to interpret.

### 12 *2.3 Land-cover model development using remote sensing and non-remote sensing predictors*

13 Land cover models were developed for each level of specificity. Seventy percent of the  
14 samples (N=1,576) were used for model calibration and 30% of the samples (N=676) were used  
15 for model validation.

16 Machine learning was used to omit redundant predictors and determine the feasibility of  
17 using the remaining predictors to predict each land cover type, given the large number of  
18 predictors and possible inter-correlations. Machine learning techniques lead to stable results  
19 when the number of predictors is large and are less affected by non-linearity and  
20 multicollinearity than other automated fitting routines (Binder and Tutz, 2008). The Breiman's  
21 random forest algorithm (Breiman, 2001) available in the “randomForest” package in R was  
22 selected in particular, because it is less susceptible to over-fitting and yields higher prediction  
23 accuracy than other machine learning algorithms (Fernández-Delgado et al., 2014). The Random

1 Forest (RF) algorithm yields an ensemble model, bagged from multiple and independent decision  
2 trees consisting of various combinations of predictors and sample subsets. The performance of  
3 the ensemble was measured with a pseudo coefficient of determination (pseudo- $R^2$ ), which is one  
4 minus the ratio of the cross-validated mean squared error (MSE) of the prediction to the variance  
5 of the observed data. As MSE or the average error between predicted and observed estimates  
6 approaches zero,  $R^2$  approaches one (perfect correlation). The importance of each predictor in  
7 the ensemble is also quantified and is defined by the percent increase in cross-validated MSE  
8 when a predictor is removed from the ensemble. Once the predictors were ranked, the “rfcv”  
9 function was used to determine the number of predictors to use to develop functional  
10 relationships for each land cover class. Rfcv computes the cross-validated MSE versus the  
11 number of predictors included in the ensemble in descending order of importance.

12 The drawback of RF is that it results in complex relationships that are difficult to  
13 interpret. Generalized Additive Models (GAMs) (Hastie and Tibshirani, 1990) were used to  
14 build functional relationships on the subsets of important predictors identified with RF, because a  
15 number of studies have successfully estimated the proportion of crop area with socio-ecological  
16 predictors and GAMs (Grace et al., 2014; Husak et al., 2008; Marshall et al., 2011); like RF,  
17 GAMs are not severely impacted by non-linear data; and unlike RF, GAMs are relatively simple  
18 and easy to interpret. Since the response variable (proportion of land cover type) was continuous  
19 and bounded from 0-100%, the data was fitted using a quasi-binomial distribution  
20 (link=logistic). The logistic GAM predicts the log-likelihood of an event (probability of  
21 success/probability of failure) using, in our case, a series of cubic spline functions:

$$\log\left(\frac{p_j}{1 - p_j}\right) = \beta_0 + \sum f_{i,j}(x_{i,j}) \quad (4)$$

1 Where  $\mathbf{p}$  is the probability of a LULC type for sample area frame  $\mathbf{j}$ ,  $\beta_0$  is the intercept, and  
2  $f_{i,j}(x_{i,j})$  is the cubic spline function for predictor  $x_i$  at sample area frame  $\mathbf{j}$ . The GAMs were  
3 developed with the “gam” function in the “mgcv” package in R. Model calibration was  
4 evaluated with explained part and overall deviance. Deviance is the log-likelihood (probability  
5 space) alternative to variance. Part deviance is the deviance explained when the target predictor  
6 is removed from a GAM minus the overall deviance. Another pseudo- $R^2$  statistic (1-model  
7 deviance/null deviance) was also computed to compare calibration with validation.

8 In order to demonstrate how the models can be used for macro-scale application, the final  
9 GAMs developed were employed to reconstruct the annual change in agriculture and natural  
10 vegetation and to perform a trend analysis from 1983-2012 at each sample area frame. Trends  
11 were estimated using the Theil-Sen technique, which computes the median of all possible  
12 pairwise slopes in a time series. The approach has been used, for example, to measure long-term  
13 trends in NDVI (de Beurs and Henebry, 2005), because it is not significantly impacted by  
14 outliers or non-linearity. The significance of each trend was assessed using the Mann-Kendall  
15 statistic. Trends were masked at the 99.9% confidence band to mask statistically insignificant  
16 trends.

### 17 **3. Results**

#### 18 *3.1 Land cover sample area frame summary*

19 The distribution of land cover over the sample area frames is illustrated with a boxplot in  
20 **Figure 3**. Agriculture and natural vegetation land cover (level one) were normally distributed,  
21 with agriculture having a higher median (54.04%) and lower spread (29.32% and 76.33% at the  
22 first and third quartiles) than natural vegetation (median=39.72%, first quartile=16.21%, and  
23 third quartile=65.67%). The proportion of urban and miscellaneous land cover was considerably

1 lower (median=4.00% and 0%, respectively) and non-linear, each having several high proportion  
2 outliers. The disaggregated land cover (level two) distributions, with the exception of crops,  
3 were non-linear with long right tails. Crops represented the largest proportion of land cover  
4 (median=37.52%) and had the largest spread (19.43% and 58.46% at the first and third quartiles),  
5 followed by savanna (median=15.79%, first quartile=3.80%, and third quartile=31.91%).  
6 Wetlands represented the smallest proportion of land cover (median=0%), with sample area  
7 frames not exceeding 75%, while forest represented the second smallest proportion of land cover  
8 (median=2.22%), but had a large number of outliers.

### 9 *3.2 Data reduction*

10 The top remote sensing and non-remote sensing predictors considered are ranked in  
11 descending order of importance for agriculture and natural vegetation using bar graphs in **Figure**  
12 **4**. The RF ensemble models using non-remote sensing predictors performed moderately well for  
13 agriculture (pseudo- $R^2=0.69$ ) and natural vegetation (pseudo- $R^2=0.69$ ), but poorly for the more  
14 non-linear distributions (urban pseudo- $R^2=0.37$  and miscellaneous pseudo- $R^2=0.50$ ). The RF  
15 ensemble models using remote sensing predictors all performed poorly: agriculture (pseudo-  
16  $R^2=0.49$ ), natural vegetation (pseudo- $R^2=0.50$ ), urban (pseudo- $R^2=0.22$ ), and miscellaneous  
17 (pseudo- $R^2=0.33$ ). It should be noted in each case however, the highest ranked remote sensing  
18 predictors resulted in lower model error than the highest ranked non-remote sensing predictors.  
19 The non-remote sensing predictors were more numerous and generated larger incremental  
20 improvements that contributed to overall greater predictive power. For the non-remote sensing  
21 ensembles, dynamic predictors were more important than slowly-changing predictors, and  
22 population density and climate predictors consistently outranked topographic or hydrologic  
23 predictors. Popd.d, popd, bio7.d, bio14.d, and bio3.d were consistently ranked the most



1 important predictors of agriculture and natural vegetation proportions. Omitting popd.d, the  
2 most important predictor for agriculture, for example, led to a more than 65% increase in  
3 ensemble MSE. Given that popd.d and popd were both important, model results were compared  
4 with popd.d and popd individually and combined as anomalies (popd.d /popd). Ensemble  
5 performance was better when the two predictors were considered separately. The most important  
6 remote sensing predictors were less influential than popd.d. strn, ampn.d, and ampl.d were more  
7 equally important for agriculture and natural vegetation, followed by phsl and phsn.

8         The importance of predictors of level two (crops, savanna, and forest) proportions are  
9 ranked in **Figure 5**. The ranking was more variable for level two classifications, but popd.d  
10 remained the most important predictor in each case. The level two RF ensemble models  
11 predicted less variability than the level one RF ensemble models and the non-remote sensing  
12 predictors outperformed the remote sensing predictors when more than the highest ranked  
13 predictors were introduced. The non-remote sensing models performed moderately well for  
14 crops (pseudo- $R^2=0.63$ ), savanna (pseudo- $R^2=0.62$ ), and forest (pseudo- $R^2=0.61$ ), but poorly for  
15 fallow (pseudo- $R^2=0.42$ ), shrubs (pseudo- $R^2=0.54$ ), wetlands (pseudo- $R^2=0.10$ ), and agroforestry  
16 (pseudo- $R^2=0.55$ ). Precipitation-based climatic predictors (bio12.d, bio13.d, bio14.d, and  
17 bio16.d) were more important in the savanna ensemble than temperature-based climatic  
18 determinants driving the crop ensemble. For the forest simulation, topographic predictors (slp  
19 and topind) were more important than most of the climatic predictors. The remote sensing  
20 ensembles performed poorly for all of the level two land cover classes: crops (pseudo- $R^2=0.46$ ),  
21 fallow (pseudo- $R^2=0.33$ ), shrubs (pseudo- $R^2=0.44$ ), savanna (pseudo- $R^2=0.44$ ), wetlands  
22 (pseudo- $R^2<1\%$ ), forest (pseudo- $R^2=0.46$ ), and agroforestry (pseudo- $R^2=0.41$ ). For crops, strn  
23 and ampl.d remained the most important predictors. Maximum annual NDVI, as captured by

1 ampl.d and ampn.d, was much more important for predicting the proportion of savanna. Unlike  
2 other ensembles, which were driven by dynamic predictors, the most important remote sensing  
3 predictors for forest cover were long-term averages.

### 4 *3.3 Building functional relationships*

5 The GAMs were developed for moderately performing land cover classes and used  
6 considerably fewer predictors than the RF ensembles, because most of the predictors in the  
7 ensembles explained very little, if any, variability. This is illustrated in **Figure 6**, which shows  
8 MSE versus the number of predictors used in the non-remote sensing and remote sensing  
9 ensembles for forest. For the non-remote sensing ensemble, MSE increased from 119.76 to  
10 120.49 after the 10<sup>th</sup> predictor and leveled off after the 13<sup>th</sup> predictor were introduced. For the  
11 remote sensing ensemble, MSE increased from 120.49 to 163.34 and levelled off after the 7<sup>th</sup>  
12 predictor was introduced. For this reason, the GAMs were built with 10-13 of the highest ranked  
13 non-remote sensing predictors and additional predictors, namely popd, were removed after  
14 redundancies were identified in the GAM component functional plots and with significance tests  
15 (not shown). GAMs were not constructed using the remote sensing predictors, because of the  
16 poor results of the ensembles and the inability of additional predictors to substantially improve  
17 the accuracy of the GAMs. Similarly, non-remote sensing GAMs were not developed for urban,  
18 miscellaneous, fallow, shurbs, or wetlands.

19 **Figures 7 and 8** show the functional relationships of the predictors used for estimating  
20 the proportion of agriculture and natural vegetation. Each model explained 61.5% (pseudo-  
21  $R^2=0.66$ ) and 61.4% (pseudo- $R^2=0.66$ ) of model deviance with nine and seven predictors,  
22 respectively. The error bars tended to be wider at proportion extremes, because fewer data  
23 points were available to train the models. The relative importance of each predictor, as defined

1 by part deviance and other calibration statistics are shown in **Table 3** for the land cover types  
2 that were considered feasible for model-building. Popd.d remained the most important predictor  
3 and uniquely explained 7.0-26.2% of model deviance. The log-likelihood of agriculture (natural  
4 vegetation) increased (decreased) rapidly as population density increased from 0 to 550  
5 people•km<sup>-2</sup>, more gradually between 550 and 1200 people•km<sup>-2</sup>, and reversed beyond 1200  
6 people•km<sup>-2</sup>. The predictive power of the topographic and climatic variables dropped off sharply  
7 after popd.d. For agriculture, bio14.d and topind were the second and third most important  
8 predictors, but explained only 1.9% and 1.6% unique deviance. As seen in the partial functional  
9 plots, the proportion of agriculture was highest in high production zones (medium population  
10 density) on ridges and crests where topind was low and for very wet tropical areas where bio14.d  
11 was high and semi-arid areas where bio14.d was low. For natural vegetation, temperature  
12 predictors, bio4.d and bio7.d, explained the second and third highest unique deviance after  
13 popd.d (2.0% and 1.3%). As seen in the functional plots, low populated areas with more  
14 temperature seasonality, or inter-annual variation, and lower bio3.d (isothermality) tended to  
15 have higher proportions of natural vegetation (savanna and shrubs). Isothermality is the ratio of  
16 mean diurnal temperature range (bio2.d) to the temperature annual range (bio7.d), which is the  
17 difference between the annual maximum and minimum temperatures. Areas that are less  
18 isothermal essentially have more pronounced seasons and are climatically less tropical. For the  
19 level two classifications, calibration was more difficult and yielded poorer relationships. Popd.d  
20 was the most important predictor and explained 7.0 – 16.4% unique deviance. The predictive  
21 power of the topographic and climatic variables was more equally distributed than for the level  
22 one classification.

1 In all cases, the  $R^2$  for the validation subset was lower than the pseudo- $R^2$  from the  
2 calibration subset: agriculture ( $\Delta R^2=-0.04$ ), natural vegetation ( $\Delta R^2=-0.01$ ), crops ( $\Delta R^2=-0.03$ ),  
3 savanna ( $\Delta R^2=-0.01$ ), and forest ( $\Delta R^2=-0.06$ ) (**Figure 9**). With the exception of the crops GAM,  
4 level two GAMs tended to under-predict high proportions of land cover (savanna and forest) and  
5 contained numerous outliers.

### 6 *3.4 Trend analysis*

7 The GAMs for agriculture and natural vegetation were used to simulate trends in the  
8 annual proportions for the sample area frames from 1983-2012 as part of the evaluation to  
9 demonstrate how the approach could be used for a retrospective analysis. The proportion of  
10 agriculture for 1983 and 2012 are shown in **Figure 10a and 10b**, while trends over the 30-year  
11 period are shown in **Figure 10c**. The high potential agricultural zone (wet highlands) in Western  
12 Kenya experienced the largest increase in simulated agricultural cover ( $> 1\%$  per year or 30%  
13 over the 30-year period). A time series of the strongest trend (1.68% per year) is shown in  
14 **Figure 10d**. Simulated population density was at 145 people $\cdot$ km $^{-2}$  in 1983 for this sample area  
15 frame, which steadily increased to 478 people $\cdot$ km $^{-2}$  in 2012. Closer to the lake, which consists  
16 of drier marginal mixed farming, trends were insignificant at the 99.9% confidence band or  
17 relatively weak ( $< 1\%$  per year). Similar patterns were seen for the marginal mixed farming and  
18 high potential agricultural zones of central Kenya. The only decrease in agricultural lands was in  
19 Kitale town (-1.40% per year). The time series is also shown in **Figure 10d**. Population growth  
20 in Kitale was 1,110 people $\cdot$ km $^{-2}$  in 1983, which is near the threshold of declining agriculture  
21 cover versus population density at 1,200 people $\cdot$ km $^{-2}$ . By 2009, when the largest decrease in  
22 agriculture cover occurred, from 51.0 to 29.5%, population density had steadily increased and  
23 surpassed another apparent threshold above 3,000 people $\cdot$ km $^{-2}$ . The direction and relative

1 magnitude of trends in natural vegetation (not shown) generally corresponded inversely to trends  
2 in agriculture, but were negatively-weak (maximum=-0.4% per year or -12% over the 30-year  
3 period).

#### 4 **5. Discussion**

5         The results make three important contributions that the land surface modeling community  
6 should consider to improve LULCC detection, particularly for SSA: 1) a socioeconomic variable  
7 (population density) was the highest ranked predictor of LULCC and had considerably more  
8 predictive power than biophysical predictors; 2) non-remote sensing predictors outperformed  
9 remote sensing predictors due to their number and the incremental improvement in the predictive  
10 power of each; and 3) coarse resolution data was able to capture general classification  
11 descriptors, but was unable to capture more detailed descriptors.

12         The global increase in agricultural land cover has been attributed to the demand for food  
13 and other agricultural commodities by a growing population (Pongratz et al., 2008). In SSA,  
14 smallholder farms, which support the majority of the labor force, are small (half are < 1.5 ha) and  
15 concentrated in densely populated areas, while large portions of arable farmland in sparsely  
16 populated areas remain underutilized (Jayne et al., 2003). This underutilization is due primarily  
17 to a lack of investment in infrastructure and unequitable tenure systems, which forces farmers to  
18 grow more on less land. This relationship is confirmed by rural population survey data in Kenya,  
19 which showed that fertilizer input use and net farm income per hectare increase until  
20 approximately 550 persons•km<sup>-2</sup> and then sharply decline, because farm sizes shrink, surplus  
21 production decreases, and farmers must adopt costlier strategies (e.g. zero-grazing) to maximize  
22 revenue (Jayne and Muyanga, 2012). The functional relationship for population density and  
23 steady increase in area under cultivation in high production zones demonstrated by the trend

1 analysis in this study, corresponds to this finding, as area under cultivation increased rapidly to  
2 approximately 550 persons•km<sup>-2</sup> and then increased more gradually with higher population  
3 density until 1200 persons•km<sup>-2</sup>. Few sample area frames had population densities greater than  
4 1,200 persons•km<sup>-2</sup>, as in Kitale town, so it is difficult to know if this functional relationship  
5 holds for very high population densities. At least to 2008, Kitale experienced a growth rate of  
6 12%, well above the national average (7%), due to persistent drought and out-migration from  
7 neighboring high production zones (Majale, 2008). Although the functional relationship for  
8 population density corroborates household surveys in Kenya and other agrarian countries in  
9 SSA, it should be further scrutinized, because land tenure in SSA is complex (Place, 2009), the  
10 dependency of LULCC predictors on location and spatial scale can be high (Rindfuss et al.,  
11 2004), and the transition from agrarian to industrialized nations may make 50-100 year  
12 projections for SSA obsolete.

13         The proposed methodology when applied to other regions of the world will undoubtedly  
14 result in a different combination of socio-ecological predictors and functional relationships,  
15 because access to land varies across agrarian and non-agrarian societies, so further study is  
16 required with observed data to develop region-specific models and validate the results for  
17 countries in SSA. Kumar et al., 2013, for example, showed that in the United States pre-1900  
18 when the country was largely agrarian and transportation networks were weak, population  
19 density and crop area were highly correlated, because crops needed to be grown close to markets.  
20 However, as the country became more industrialized and transportation networks improved,  
21 farmers moved to more biophysically suitable areas away from city centers, making biophysical  
22 determinants of crop area more important than population density in the latter half of the 20<sup>th</sup>  
23 century. Whether the analyses are performed in agrarian or non-agrarian regions, extensive

1 preparation of observation data will be required, because the data used in this study, namely  
2 consistent sample area frames at a spatial resolution appropriate for land surface modeling and  
3 spanning multiple climatic zones through time, is quite unique.

4 Population density estimates vary widely (Wilson, 2014) and given its fundamental  
5 importance to the proposed model framework, future work should aim to integrate a more  
6 dynamic product that better accounts for inter-annual variability and realistic representation of  
7 current and projected population density. To the authors' knowledge, this was the first attempt  
8 to make a population product dynamic. However, the approach is essentially tracking decadal  
9 trends that explain a significant portion of inter-annual variability. In reality, population density  
10 can show high inter-annual variability due to migration and other factors. Regarding the product  
11 itself, changes in population density do not necessarily "grow" from transportation networks and  
12 are influenced by important feedbacks now and in the future. In addition, the extrapolation  
13 method used is efficient and can be projected indefinitely, but does not capture complex  
14 demographics that other methods do and can lead to "runaway" growth/decline and unrealistic  
15 mid- to late- 21<sup>st</sup> century projections for scenario-building (Baker et al., 2008). Finally, there is  
16 no consensus on which population product to use however, in the future, other products (e.g.  
17 Afripop) should be compared against the product used here, used to adjust growth/decay  
18 coefficient for population density estimates beyond 2000, or combined to make a model  
19 ensemble.

20 This paper highlights the importance of gridded socioeconomic data in mapping LULCC,  
21 but gridded macro-scale datasets are almost exclusively biophysical in nature. The biggest gains  
22 in LULCC prediction could be made, therefore, by developing gridded macro-scale  
23 socioeconomic data from existing country-level products, such as the Human Development

1 Index. More minor gains could be made by integrating biophysical predictors not used in this  
2 study, such as soil type and properties. Gridded soils data exists globally from the International  
3 Soil Reference and Information Center, but was not considered in this study, because it is a one-  
4 time value and does not capture the dynamic nature of soils or its complex relationship with  
5 LULCC. A dynamic soils product was recently developed for the MODIS era (see Vågen et al.,  
6 2016) and could be a powerful tool for LULCC detection, especially if it is back-casted over the  
7 full temporal range of other predictors. Many biophysical predictors are available mid- and late-  
8 21<sup>st</sup> century and are therefore widely used for prospective analyses, so methods should be  
9 explored to project soils and socio-economic data into the future to improve LULCC estimates.

10 Grace et al. (2014) developed GAMs to predict cropped area in Kenya using biophysical  
11 predictors (rainfall, elevation, NDVI, slope, and the topographic wetness index) and explained  
12 much of the deviance in cropped area (41.9-81.4%). Although the models used different  
13 predictors for different years and production zones, and the definition of cropped area and the  
14 degree of functional smoothing were not explicit, the study highlights that the inter-correlation  
15 among predictors may be obscuring the importance of biophysical determinants. Specifically,  
16 population density tends to be highly correlated with and could be suppressing the explanatory  
17 power of biophysical predictors, though the partial deviance statistics did not reflect this. In  
18 addition, the random forest algorithm accounts for inter-correlation to some degree, but other  
19 techniques could be introduced to further reduce these effects. For example, Principal  
20 Components Analysis could be used to develop temperature and precipitation indices that  
21 integrate all or some of the BIOCLIM predictors.

22 Phenological patterns extracted from continuous Earth observation based NDVI have  
23 been widely used to map LULCC over long time periods, given the lack of higher spatial and



1 spectral resolution data before the MODIS era (Ali et al., 2014; Bie et al., 2012). These studies  
2 show that vegetation periodicity is highly variable for a given land cover type and that long-term  
3 averages of phenological predictors are more reliable for mapping LULC. In this study, many of  
4 the important remote sensing predictors (particularly for forests) were long-term averages, but  
5 they still under-predicted LULCC when compared against non-remote sensing predictors, which  
6 were more numerous and resulted in larger incremental improvements to model accuracy.  
7 Perhaps the main difficulty in using long-term Earth observation data for LULCC estimation is  
8 the coarseness of the data and the rapid change in vegetation that often occurs over small spatial  
9 scales. Population density, which was a much stronger predictor, on the other hand, may well be  
10 captured using coarse resolution data, because this predictor changes more gradually over space.  
11 An analysis of the non-remote sensing and remote sensing predictors together revealed that for  
12 agriculture, natural vegetation, savanna, and forest cover, Earth observation data provided an  
13 additional 1-2% explained deviance. If the long-term average remote sensing predictors could be  
14 downscaled using MODIS or Landsat data and then aggregated to 5x5 km<sup>2</sup> resolution with  
15 distribution moments as predictors, for example, the explanatory power of non-remote sensing  
16 predictors could be further enhanced for retrospective analyses. Another avenue worth exploring  
17 could involve using downscaled long-term average remote sensing predictors to develop 5x5 km<sup>2</sup>  
18 probabilities as in the Pengra et al., 2015 dataset to evaluate the non-remote sensing models  
19 proposed here.

20 The evaluation of the models at two levels of specificity revealed that coarse resolution is  
21 able to better simulate general descriptors, such as natural vegetation, but is poorer at predicting  
22 more detailed descriptors, such as forest. Each of the more detailed random forest ensembles  
23 with non-remote sensing predictors had  $\Delta R^2$ s of -0.06, -0.07, and -0.07 for crop, savanna, and

1 forest over agriculture and natural vegetation, respectively. Part of this discrepancy can be  
2 attributed to the increased interpretation uncertainty, as interpreters find it more challenging to  
3 distinguish between more detailed LULC types. In addition, coarse resolution data may not be  
4 able to capture the level of heterogeneity in the area sample frames needed to distinguish land  
5 use/cover-specific socio-ecological patterns and properties.

## 6 **6. Conclusion**

7 This study developed and evaluated a simple method to provide consistent estimates of  
8 LULCC annually over 30 years at 5x5 km<sup>2</sup> resolution using non-parametric functional  
9 relationships with a small subset of socio-ecological predictors ( $p \leq 10$ ). Functional relationships  
10 were developed after data mining 43 geospatial datasets that are available seamlessly across  
11 SSA, which can be used for retrospective or prospective mid- and late-21<sup>st</sup> century analyses as  
12 well. The relationships are intuitive and tunable, making their use practical for decision-makers  
13 to identify intervention hotspots and develop land management scenarios. Model validation,  
14 performed with the proportion of major land cover types in Kenya over a 30-year period,  
15 revealed that a number of activities should be performed to improve the predictive power of the  
16 models for practical use. These activities should include integrating improved existing or newly  
17 developed geospatial (particularly socioeconomic) datasets into the proposed model framework.  
18 With these improvements, land surface and LULCC modelling could be greatly enhanced and  
19 the consequence of the latter on the earth system more fully understood. In an upcoming study,  
20 the modeling approach proposed here will be used with a newly acquired sample area frame  
21 dataset to estimate baseline LULCC and project land suitability across SSA mid-21<sup>st</sup> century  
22 with AFRICLIM and other geospatial data.

1 **Acknowledgements**

2           The work summarized in this manuscript was primarily funded through support by the  
3 CGIAR research program on Climate Change, Agriculture and Food Security for the project,  
4 titled “Multi-disciplinary species distribution modelling for “climate smart” agriculture in East  
5 Africa.” Additional support for the early field and aerial surveys was supported by the Kenyan  
6 Lake Basin Development Authority and Ministry of Planning and National Development. We  
7 would like to extend our special thanks to Sandra Nakibilango, Dorcas Ninsiima, Charles Ngugi,  
8 Patrick Ojorot and Samuel Olowo who interpreted the aerial photos taken for this project;  
9 Bernadette Apio who coordinated the interpreter team; and Juliet Kyakobyewo who entered the  
10 data into our database. Finally, we would like to thank Dr. Eike Luedeling, who initially led and  
11 organized the project.

## References

- 1 Alcamo, J., Schaldach, R., Koch, J., Kölking, C., Lapola, D., Priess, J., 2011. Evaluation of an  
2 integrated land use change model including a scenario analysis of land use change for  
3 continental Africa. *Environ. Model. Softw.* 26, 1017–1027.
- 4 Ali, A., de Bie, C.A.J.M., Skidmore, A.K., Scarrott, R.G., Lymberakis, P., 2014. Mapping the  
5 heterogeneity of natural and semi-natural landscapes. *Int. J. Appl. Earth Obs.*  
6 *Geoinformation* 26, 176–183.
- 7 Anderson-Teixeira, K.J., DeLUCIA, E.H., 2011. The greenhouse gas value of ecosystems. *Glob.*  
8 *Change Biol.* 17, 425–438.
- 9 Baker, J., Ruan, X., Alcantara, A., Jones, T., Watkins, K., McDaniel, M., Frey, M., Crouse, N.,  
10 Rajbhandari, R., Morehouse, J., Sanchez, J., Inglis, M., Baros, S., Penman, S., Morrison,  
11 S., Budge, T., Stallcup, W., 2008. Density-dependence in urban housing unit growth: An  
12 evaluation of the Pearl-Reed model for predicting housing unit stock at the census tract  
13 level. *J. Econ. Soc. Meas.* 33, 155–163.
- 14 Ban, Y., Gong, P., Giri, C., 2015. Global land cover mapping using Earth observation satellite  
15 data: Recent progresses and challenges. *ISPRS J. Photogramm. Remote Sens.* 103, 1–6.
- 16 Bie, C.A.J.M. de, Nguyen, T.T.H., Ali, A., Scarrott, R., Skidmore, A.K., 2012. LaHMa: a  
17 landscape heterogeneity mapping method using hyper-temporal datasets. *Int. J. Geogr.*  
18 *Inf. Sci.* 26, 2177–2192.
- 19 Binder, H., Tutz, G., 2007. A comparison of methods for the fitting of generalized additive  
20 models. *Stat. Comput.* 18, 87–99.
- 21 Breiman, L., 2001. Random Forests. *Mach. Learn.* 45, 5–32.

1 Carrão, H., Gonalves, P., Caetano, M., 2010. A Nonlinear Harmonic Model for Fitting Satellite  
2 Image Time Series: Analysis and Prediction of Land Cover Dynamics. *IEEE Trans.*  
3 *Geosci. Remote Sens.* 48, 1919–1930.

4 Chaney, N.W., Sheffield, J., Villarini, G., Wood, E.F., 2014. Development of a High-Resolution  
5 Gridded Daily Meteorological Dataset over Sub-Saharan Africa: Spatial Analysis of  
6 Trends in Climate Extremes. *J. Clim.* 27, 5815–5835.

7 Chen, J., Chen, J., Liao, A., Cao, X., Chen, L., Chen, X., He, C., Han, G., Peng, S., Lu, M.,  
8 Zhang, W., Tong, X., Mills, J., 2015. Global land cover mapping at 30 m resolution: A  
9 POK-based operational approach. *ISPRS J. Photogramm. Remote Sens., Global Land*  
10 *Cover Mapping and Monitoring* 103, 7–27.

11 Davin, E.L., de Noblet-Ducoudré, N., 2010. Climatic Impact of Global-Scale Deforestation:  
12 Radiative versus Nonradiative Processes. *J. Clim.* 23, 97–112.

13 Davis, H.C., 1995. *Demographic Projection Techniques for Regions and Smaller Areas: A*  
14 *Primer*. UBC Press.

15 de Beurs, K.M., Henebry, G.M., 2005. A statistical framework for the analysis of long image  
16 time series. *Int. J. Remote Sens.* 26, 1551–1573.

17 DeFries, R.S., Field, C.B., Fung, I., Justice, C.O., Los, S., Matson, P.A., Matthews, E., Mooney,  
18 H.A., Potter, C.S., Prentice, K., Sellers, P.J., Townshend, J.R.G., Tucker, C.J., Ustin,  
19 S.L., Vitousek, P.M., 1995. Mapping the land surface for global atmosphere-biosphere  
20 models: Toward continuous distributions of vegetation’s functional properties. *J.*  
21 *Geophys. Res. Atmospheres* 100, 20867–20882.

1 Deichmann, U., 1996. A Review of Spatial Population Database Design and Modeling  
2 (Technical Report No. 96-3). National Center for Geographic Information and Analysis,  
3 Santa Barbara, CA.

4 Eastman, R., Sangermano, F., Ghimire, B., Zhu, H., Chen, H., Neeti, N., Cai, Y., Machado, E.A.,  
5 Crema, S.C., 2009. Seasonal trend analysis of image time series. *Int. J. Remote Sens.* 30,  
6 2721–2726. doi:10.1080/01431160902755338

7 EcoSystems Ltd, 1987. Integrated Land Use Database for Kenya. Ministry of Planning & Natural  
8 Development, Nairobi, Kenya.

9 EcoSystems Ltd, 1983. Integrated Land Use Survey: Final Report. Lake Basin Development  
10 Authority, Kisumu, Kenya.

11 Elzhov, T.V., Mullen, K.M., Spiess, A.N., Bolker, B., 2015. Package “minpack.lm.”

12 Fernández-Delgado, M., Cernadas, E., Barro, S., Amorim, D., 2014. Do We Need Hundreds of  
13 Classifiers to Solve Real World Classification Problems? *J Mach Learn Res* 15, 3133–  
14 3181.

15 Funk, C., Peterson, P., Landsfeld, M., Pedreros, D., Verdin, J., Rowland, J., Romero, B., Husak,  
16 G.J., Michaelsen, J., Verdin, A., 2014. A Quasi-Global Precipitation Time Series for  
17 Drought Monitoring (No. 832), U.S. Geological Survey Data Series. U.S. Geological  
18 Survey, Washington, D.C.

19 Funk, C., Verdin, A., Michaelsen, J., Peterson, P., Pedreros, D., Husak, G., 2015. A global  
20 satellite assisted precipitation climatology. *Earth Syst. Sci. Data Discuss.* 8, 401–425.

21 Giri, C., Pengra, B., Long, J., Loveland, T.R., 2013. Next generation of global land cover  
22 characterization, mapping, and monitoring. *Int. J. Appl. Earth Obs. Geoinformation* 25,  
23 30–37.

1 Grace, K., Husak, G., Bogle, S., 2014. Estimating agricultural production in marginal and food  
2 insecure areas in Kenya using very high resolution remotely sensed imagery. *Appl.*  
3 *Geogr.* 55, 257–265.

4 Grace, K., Husak, G.J., Harrison, L., Pedreros, D., Michaelsen, J., 2012. Using high resolution  
5 satellite imagery to estimate cropped area in Guatemala and Haiti. *Appl. Geogr.* 32, 433–  
6 440.

7 Hansen, M.C., Stehman, S.V., Potapov, P.V., 2010. Quantification of global gross forest cover  
8 loss. *Proc. Natl. Acad. Sci.* 107, 8650–8655.

9 Hansen, M.C., Loveland, T.R., 2012. A review of large area monitoring of land cover change  
10 using Landsat data. *Remote Sens. Environ., Landsat Legacy Special Issue* 122, 66–74.

11 Hargreaves, G.H., Samani, Z.A. 1985. Reference Crop Evapotranspiration from Temperature.  
12 *Appl. Eng. Agric.* 1, 96–99.

13 Hastie, T.J., Tibshirani, R.J., 1990. *Generalized Additive Models*. CRC Press.

14 Held, I.M., Soden, B.J., 2006. Robust Responses of the Hydrological Cycle to Global Warming.  
15 *J. Clim.* 19, 5686–5699.

16 Heistermann, M., Müller, C., Ronneberger, K., 2006. Land in sight?: Achievements, deficits and  
17 potentials of continental to global scale land-use modeling. *Agric. Ecosyst. Environ.* 114,  
18 141–158.

19 Hijmans, R.J., Phillips, S., Leathwick, J., Elith, J., 2015. Package “dismo.”

20 Hijmans, R.J., S.E. Cameron, J.L. Parra, P.G. Jones and A. Jarvis, 2005. Very high resolution  
21 interpolated climate surfaces for global land areas. *International Journal of Climatology.*  
22 25, 1965-1978.

1 Husak, G.J., Marshall, M.T., Michaelsen, J., Pedreros, D., Funk, C., Galu, G., 2008. Crop area  
2 estimation using high and medium resolution satellite imagery in areas with complex  
3 topography. *J. Geophys. Res. Atmospheres* 113, D14112.

4 Jayne, T.S., Yamano, T., Weber, M.T., Tschirley, D., Benfica, R., Chapoto, A., Zulu, B., 2003.  
5 Smallholder income and land distribution in Africa: implications for poverty reduction  
6 strategies. *Food Policy* 28, 253–275.

7 Jayne, T.S., Muyanga, M., 2012. Land constraints in Kenya’s densely populated rural areas:  
8 implications for food policy and institutional reform. *Food Secur.* 4, 399–421.

9 Kumar, S., Merwade, V., Rao, P.S.C., Pijanowski, B.C., 2013. Characterizing Long-Term Land  
10 Use/Cover Change in the United States from 1850-2000 Using a Nonlinear Bi-analytical  
11 Model. *Ambio* 42, 285-297.

12 Lambin, E.F., Geist, H.J., Lepers, E., 2003. Dynamics of Land-Use and Land-Cover Change in  
13 Tropical Regions. *Annu. Rev. Environ. Resour.* 28, 205–241.

14 Lamprey, R.H., 2013. Aerial Point Sampling (APS) Survey: Lake Basin, Machakos and  
15 Makueni, Kenya, 2012-13, Nairobi, Kenya.

16 Lawrence, P. J., Feddema, J. J., Bonan, G. B., Meehl, G. A., O’Neill, B. C., Oleson, K.  
17 W., Levis, S., Lawrence, D. M., Kluzek, E., Lindsay, K., and Thornton, P. E., 2012.  
18 Simulating the Biogeochemical and Biogeophysical Impacts of Transient Land Cover  
19 Change and Wood Harvest in the Community Climate System Model (CCSM4) from  
20 1850 to 2100. *J. Climate* 25, 3071–3095.

21 Lepers, E., Lambin, E.F., Janetos, A.C., DeFries, R., Achard, F., Ramankutty, N., Scholes, R.J.,  
22 2005. A Synthesis of Information on Rapid Land-cover Change for the Period 1981–  
23 2000. *BioScience* 55, 115–124.



1 Majale, M., 2008. Employment creation through participatory urban planning and slum  
2 upgrading: The case of Kitale, Kenya. *Habitat Int., Labour in Urban Areas* 32, 270–282.

3 Makarieva, A.M., Gorshkov, V.G., Li, B.-L., 2013. Revisiting forest impact on atmospheric  
4 water vapor transport and precipitation. *Theor. Appl. Climatol.* 111, 79–96.

5 Marshall, M.T., Husak, G.J., Michaelsen, J., Funk, C., Pedreros, D., Adoum, A., 2011. Testing a  
6 high-resolution satellite interpretation technique for crop area monitoring in developing  
7 countries. *Int. J. Remote Sens.* 32, 7997–8012. doi:10.1080/01431161.2010.532168

8 Meiyappan, P., Dalton, M., O’Neill, B.C., Jain, A.K., 2014. Spatial modeling of agricultural land  
9 use change at global scale. *Ecol. Model.* 291, 152–174.

10 Moré, J.J., 1978. The Levenberg-Marquardt algorithm: Implementation and theory, in: Watson,  
11 G.A. (Ed.), *Numerical Analysis, Lecture Notes in Mathematics*. Springer Berlin  
12 Heidelberg, pp. 105–116.

13 Ngetich, K.F., Mucheru-Muna, M., Mugwe, J.N., Shisanya, C.A., Diels, J., Mugendi, D.N.,  
14 2014. Length of growing season, rainfall temporal distribution, onset and cessation dates  
15 in the Kenyan highlands. *Agric. For. Meteorol.* 188, 24–32.

16 Norton-Griffiths, M., 1988. Aerial Point Sampling for Land Use Surveys. *J. Biogeogr.* 15, 149–  
17 156.

18 Olofsson, P., Stehman, S.V., Woodcock, C.E., Sulla-Menashe, D., Sibley, A.M., Newell, J.D.,  
19 Friedl, M.A., Herold, M., 2012. A global land-cover validation data set, part I:  
20 fundamental design principles. *Int. J. Remote Sens.* 33, 5768–5788.

21 Pengra, B., Long, J., Dahal, D., Stehman, S.V., and Loveland, T.R., 2015. A global reference  
22 database from very high resolution commercial satellite data and methodology for

1 application to Landsat derived 30m continuous field tree cover data *Remote Sens.*  
2 *Environ.* 165, 234-248.

3 Pielke, R.A., Pitman, A., Niyogi, D., Mahmood, R., McAlpine, C., Hossain, F., Goldewijk, K.K.,  
4 Nair, U., Betts, R., Fall, S., Reichstein, M., Kabat, P., de Noblet, N., 2011. Land use/land  
5 cover changes and climate: modeling analysis and observational evidence. Wiley  
6 *Interdiscip. Rev. Clim. Change* 2, 828–850.

7 Pinzon, J.E., Tucker, C.J., 2014. A Non-Stationary 1981–2012 AVHRR NDVI3g Time Series.  
8 *Remote Sens.* 6, 6929–6960.

9 Pitman, A.J., 2003. The evolution of, and revolution in, land surface schemes designed for  
10 climate models. *Int. J. Climatol.* 23, 479–510.

11 Place, F., 2009. Land Tenure and Agricultural Productivity in Africa: A Comparative Analysis of  
12 the Economics Literature and Recent Policy Strategies and Reforms. *World Dev.*, The  
13 Limits of State-Led Land Reform 37, 1326–1336.

14 Platts, P.J., Omeny, P.A., Marchant, R., 2014. AFRICLIM: high-resolution climate projections  
15 for ecological applications in Africa. *Afr. J. Ecol.* 53, 103–108.

16 Pongratz, J., Reick, C., Raddatz, T., Claussen, M., 2008. A reconstruction of global agricultural  
17 areas and land cover for the last millennium. *Glob. Biogeochem. Cycles* 22, GB3018.

18 Pricope, N.G., Husak, G., Lopez-Carr, D., Funk, C., Michaelsen, J., 2013. The climate-  
19 population nexus in the East African Horn: Emerging degradation trends in rangeland and  
20 pastoral livelihood zones. *Glob. Environ. Change* 23, 1525–1541.

21 Rindfuss, R.R., Walsh, S.J., Turner, B.L., Fox, J., Mishra, V., 2004. Developing a science of land  
22 change: Challenges and methodological issues. *Proc. Natl. Acad. Sci. U. S. A.* 101,  
23 13976–13981.

1 Rounsevell, M.D.A., Arneth, A., Alexander, P., Brown, D.G., de Noblet-Ducoudré, N., Ellis, E.,  
2 Finnigan, J., Galvin, K., Grigg, N., Harman, I., Lennox, J., Magliocca, N., Parker, D.,  
3 O'Neill, B.C., Verburg, P.H., Young, O., 2014. Towards decision-based global land use  
4 models for improved understanding of the Earth system. *Earth Syst Dynam* 5, 117–137.

5 Schaldach, R., Priess, J.A., 2008. Integrated Models of the Land System: A Review of Modelling  
6 Approaches on the Regional to Global Scale. *Living Rev. Landsc. Res.* 2.  
7 doi:10.12942/lrlr-2008-1

8 Sheffield, J., Goteti, G., Wood, E.F., 2006. Development of a 50-Year High-Resolution Global  
9 Dataset of Meteorological Forcings for Land Surface Modeling. *J. Clim.* 19, 3088–3111.

10 Shevliakova, E., Pacala, S. W., Malyshev, S., Hurtt, G. C., Milly, P. C. D., Caspersen,  
11 J. P., Sentman, L. T., Fisk, J. P., Wirth, C., and Crevoisier, C., 2009. Carbon cycling under  
12 300 years of land use change: Importance of the secondary vegetation sink. *Biogeochem.*  
13 *Cy.* 23, GB2022, doi:10.1029/2007gb003176.

14 Sterling, S.M., Ducharne, A., Polcher, J., 2013. The impact of global land-cover change on the  
15 terrestrial water cycle. *Nat. Clim. Change* 3, 385–390.

16 Tian, F., Fensholt, R., Verbesselt, J., Grogan, K., Horion, S., Wang, Y., 2015. Evaluating  
17 temporal consistency of long-term global NDVI datasets for trend analysis. *Remote Sens.*  
18 *Environ.* 163, 326–340.

19 Turner, B.L., Janetos, A.C., Verbug, P.H., Murray, A.T., 2013. Land System Architecture: Using  
20 Land Systems to Adapt and Mitigate Global Environmental Change. *Glob. Environ.*  
21 *Change* 232395-397.

- 1 Turner, B.L., Lambin, E.F., Reenberg, A., 2007. The emergence of land change science for  
2 global environmental change and sustainability. *Proc. Natl. Acad. Sci.* 104, 20666–  
3 20671.
- 4 UNEP, 2008. *Africa: Atlas of Our Changing Environment*. UN Environment Programme,  
5 Nairobi, Kenya.
- 6 Vågen, T.-G., Winowiecki, L.A., Tondoh, J.E., Desta, L.T., Gumbricht, T., 2016. Mapping of  
7 soil properties and land degradation risk in Africa using MODIS reflectance. *Geoderma*  
8 263, 216–225.
- 9 van Asselen, S., Verburg, P.H., 2013. Land cover change or land-use intensification: simulating  
10 land system change with a global-scale land change model. *Glob. Change Biol.* 19, 3648–  
11 3667.
- 12 Veldkamp, A., Fresco, L.O., 1996. CLUE-CR: An integrated multi-scale model to simulate land  
13 use change scenarios in Costa Rica. *Ecol. Model.* 91, 231–248.
- 14 Verburg, P.H., Neumann, K., Nol, L., 2011. Challenges in using land use and land cover data for  
15 global change studies. *Glob. Change Biol.* 17, 974–989
- 16 Verburg, P.H., Soepboer, W., Veldkamp, A., Limpiada, R., Espaldon, V., Mastura, S.S.A., 2002.  
17 Modeling the Spatial Dynamics of Regional Land Use: The CLUE-S Model. *Environ.*  
18 *Manage.* 30, 391–405.
- 19 Ward, D. S., Mahowald, N. M., and Kloster, S., 2014. Potential climate forcing of land use and  
20 land cover change. *Atmos. Chem. Phys.* 14, 12701–12724.
- 21 Wilson, T., 2014. New Evaluations of Simple Models for Small Area Population Forecasts:  
22 Small Area Population Forecasts. *Popul. Space Place* 21, 335–353.

1 Yu, L., Liang, L., Wang, J., Zhao, Y., Cheng, Q., Hu, L., Liu, S., Yu, L., Wang, X., Zhu, P., Li,  
2 X., Xu, Y., Li, C., Fu, W., Li, X., Li, W., Liu, C., Cong, N., Zhang, H., Sun, F., Bi, X.,  
3 Xin, Q., Li, D., Yan, D., Zhu, Z., Goodchild, M.F., Gong, P., 2014. Meta-discoveries  
4 from a synthesis of satellite-based land-cover mapping research. *Int. J. Remote Sens.* 35,  
5 4573–4588.