1	Pixels, blocks of pixels, and polygons: choosing a spatial unit for thematic accuracy assessment
2	
3	Stephen V. Stehman <sup>1</sup> and James D. Wickham <sup>2</sup>
4	
5	<sup>1</sup> State University of New York, College of Environmental Science and Forestry, 1 Forestry Drive,
6	Syracuse, NY 13210 USA
7	
8	<sup>2</sup> U.S. Environmental Protection Agency, Environmental Sciences Division, Research Triangle Park, NC
9	27711 USA
10	
11	
12	* Corresponding author:
13	Stephen Stehman
14	SUNY College of Environmental Science and Forestry
15	1 Forestry Drive
16	Syracuse, NY 13210 USA
17	Email: <u>svstehma@syr.edu</u>
18	Voice: 1 315 470 6692
19	Fax: 1 315 470 6535
20	
21	
22	

#### 23 Abstract

Pixels, blocks of pixels, and polygons are all potentially viable spatial assessment units for 24 25 conducting an accuracy assessment. We develop a population-based statistical framework to examine 26 how the spatial unit chosen affects the outcome of an accuracy assessment. The population is 27 conceptualized as a difference map created by overlaying a complete coverage reference classification 28 and the target map being evaluated. The per-class areas of agreement and disagreement derived from this 29 population are summarized by a population error matrix and accuracy parameters (e.g., overall, user's and 30 producer's accuracies). The population and values of the accuracy parameters are strongly affected by the 31 protocols implemented for the response design which include the choice of spatial unit, how within-unit 32 homogeneity is addressed when assigning class labels, and the definition of agreement between the 33 reference and map classification. Several complete coverage populations are used to illustrate how 34 accuracy results are affected by the spatial unit chosen for the assessment and also to evaluate how spatial 35 misregistration of the map and reference locations impacts accuracy results for different spatial units. The 36 sampling design implemented for accuracy assessment does not change the population or values of the 37 accuracy parameters, but the choice of spatial unit will influence decisions regarding use of strata and 38 clusters in the design. A universally best spatial assessment unit does not exist, so it is critical to 39 recognize how the population, values of the accuracy parameters, and sampling design are impacted by 40 the choice of spatial unit.

41

Key Words: response design; stratified sampling; cluster sampling; land cover; change; location error;
43

#### 44 **1. Introduction**

Accuracy assessment is an established component of the process of creating and distributing
thematic maps. The fundamental basis of an accuracy assessment is a location-specific comparison
between the map classification and the ground condition or "reference" classification. Reporting thematic
map accuracy in the form of an error matrix is a standard practice when the map and reference

49 classifications are based on categories such as land-cover classes (Story and Congalton 1986; Congalton 50 and Green 1999, 2008). However, opinions vary on the appropriate spatial unit for comparing the map 51 and reference classifications to obtain the data summarized by an error matrix. Pixels, blocks of pixels 52 (i.e., square arrays of pixels), and polygons are the spatial units commonly used. The lack of consensus 53 regarding choice of assessment unit is evident from the 33 map accuracy assessments reviewed by 54 Stehman and Czaplewski (1998), who reported that 14 assessments used a pixel as the spatial unit, 10 55 assessments used a square block of pixels (e.g., 2x2, 3x3), and 9 assessments used a polygon. Strahler et 56 al. (2006), Janssen and van der Wel (1994) and Richards (1996) support a pixel-based assessment, 57 whereas Congalton and Green (1999) recommend using a block of pixels or a polygon. 58 Because it is impractical to obtain a census of the reference classification, a sample of units is 59 selected from the region of interest (ROI) and the reference classification is obtained for each sample unit. 60 The response design is the protocol for determining the reference classification of a sampled assessment 61 unit (Stehman and Czaplewski 1998). Key decisions for the response design include choosing the spatial 62 assessment unit, specifying how the reference information will be obtained (e.g., from field observation or 63 high resolution imagery) and stating how agreement between the map and reference classification will be 64 defined. The results of an accuracy assessment are strongly influenced by the choice of spatial unit and 65 definition of agreement. The choice of sampling design does not affect the values of the accuracy 66 parameters, but different sampling designs will differ in terms of the precision of the estimates of the 67 accuracy parameter values.

The objective of this article is to elucidate how the choice of the spatial unit for accuracy assessment affects the implementation and results of the assessment. We develop a population-based conceptual framework to demonstrate how the population and the parameter values targeted by an accuracy assessment are determined by the choice of spatial unit and response design (Section 2). Several example populations provide the basis for illustrating numerically how different choices for the spatial unit and response design translate to different values of the accuracy parameters and therefore different results. The same population framework is then applied to illustrate how the spatial assessment unit impacts the

outcome of an accuracy assessment when location error is present (i.e., when the spatial units of the map classification and the spatial units of the reference classifications are not aligned). The practical ramifications of the choice of the spatial unit on the sampling design are discussed in Section 3. The key concepts and results developed in the article are previewed in Table 1.

79

#### 80 2. A Population Framework for Accuracy Assessment

81 A statistical population is defined as the collection of all elements of interest and one or more 82 quantities ("variables of study") associated with each element (Särndal et al. 1992, p.5). A parameter is a 83 function of the quantities assigned to each element where this function incorporates all elements of the 84 population. For example, population means, totals, and ratios (e.g., a ratio of means of two variables) are 85 common parameters. For accuracy assessment, a population can be defined as all spatial units forming a 86 partition of the region of interest (ROI) where the spatial unit could be a pixel, a square block of pixels, or 87 a polygon, and the observations or variables of study associated with each spatial unit are obtained from a 88 complete coverage reference classification (i.e., the reference map) and the map to be evaluated (i.e., the 89 target map). The reference map is itself a population in which the observation on each spatial unit is the 90 reference class associated with that unit. Similarly, the map being evaluated may be viewed as a 91 population with the map class assigned to each unit being the observation of interest. The reference 92 population represents the true condition on the ground and the map population represents the classified or 93 predicted condition. The difference between these two populations is the population of interest for 94 accuracy assessment, where an observation of this population could be an indicator variable representing 95 whether a spatial unit is classified correctly (i.e., the variable is assigned the value of 1 if the reference 96 label and map label agree and the value of 0 if the labels disagree) or a quantity representing the degree of 97 agreement between the map and reference classification for that spatial unit. In the case of a binary 98 representation of "agree" or "disagree", the population can be created by overlaying the reference and 99 target maps and determining the per-class area of agreement and the area of disagreement by type of 100 misclassification. It is these class-specific areas of agreement and disagreement that are summarized by

the population error matrix commonly used to describe accuracy. Our evaluation will focus on the
 parameters overall, user's, and producer's accuracies computed from a population error matrix.

103 The perspective underlying this population framework is developed from the design-based 104 inference framework (Särndal et al. 1992, sections 2.5 and 2.7). In this framework, the observation or 105 measurement taken on each unit of the population is regarded as a fixed constant, not a random variable, 106 and uncertainty is attributable to the randomization distribution resulting from the sampling design (i.e., 107 uncertainty is represented by the variation of the estimate of the parameter of interest over the set of all 108 possible samples). The value of a parameter is computed from a census of the population.

In practice, a reference map does not exist. Consequently, the reference classification will be available for only a sample of the ROI, and the error matrix and accuracy parameters must be estimated from this sample. The estimation objective still targets the error matrix and associated parameters of the population created by overlaying complete coverage reference and target maps. Other parameters that quantitatively compare complete coverage reference and target maps may also be of interest (e.g., Power et al. 2001, Dungan 2006, Hargrove et al. 2006, and White 2006) but will not be discussed further.

115 The population described above is the foundation of an area-based accuracy assessment in which 116 the objective is to characterize accuracy in terms of area or proportion of area correctly classified and area 117 or proportion of area misclassified. An alternate approach leading to a different population is a per-118 polygon or per-object accuracy assessment in which the focus is on whether the polygons or objects 119 mapped are classified correctly as individual entities. For example, Smith et al. (2002) investigated the 120 distribution of small water bodies because of their potential significance for altering sediment delivery, 121 and an important feature of this distribution is the number of small water bodies by geographic region. In 122 this example, the accuracy assessment could focus on the small water bodies as countable objects (i.e., 123 distinct individual entities) as opposed to quantifying the area of small water bodies. In such a per-object 124 assessment (in which the objects are defined by the map classification), an object may be considered 125 correctly labeled if the majority of the object's area, as determined from the reference classification, 126 corresponds to the map label. A per-object assessment thus may employ a fundamentally different

definition of agreement than an area-based assessment and the different approaches may produce differentaccuracy outcomes. The focus of this article will be area-based assessment.

129 The three spatial units we discuss are a pixel, a block of pixels, and a polygon. We assume that 130 each pixel is comprised entirely of a single class, but a block is not similarly assumed to be internally 131 homogeneous and could be comprised of a mix of classes (the "mixed pixel" case would be treated in the 132 same manner as a heterogeneous block). Of the three choices of spatial assessment unit, polygons are 133 commonly viewed as more natural on the basis that they represent real features of the landscape. This 134 viewpoint is appropriate if the polygons are defined by the reference classification, but not if polygons are 135 defined by the map classification. The polygons often used in practice for accuracy assessment are 136 defined by the map classification, and these map polygons are not always features of interest. For 137 example, map polygons may be added, eliminated, or changed as the map is revised prior to completion 138 of the final map product. Congalton and Green (2008, figure 5.5) present an illuminating example of this 139 phenomenon in which a sample map polygon selected on the basis of a preliminary version of a map turns 140 out to include portions of the area of three different polygons of the final map. Consequently, the polygon 141 assessment unit in Congalton and Green's example is not a real feature of the final map and such a 142 polygon would not be relevant to the objective of assessing the accuracy of the final map. Among the 143 options for spatial assessment unit, polygons defined by the reference classification are closest to 144 representing actual earth surface features. However, it is difficult in practice to design and implement an 145 accuracy assessment using reference polygons as the spatial assessment units (see Section 3).

The minimum mapping unit (mmu) specified for the target map may be taken into consideration when choosing the spatial unit. It is not sufficient to state that the mmu is the spatial unit for the assessment because simply specifying the mmu does not unambiguously partition the ROI. One option for taking into account the mmu would be to partition the ROI into square blocks where the area of each block is equivalent to the area defined by the mmu (e.g., a 9-pixel mmu would translate to a 3x3 pixel block unit). However, as is generally true of any partition formed by blocks, some blocks will contain a mix of classes and this within-block heterogeneity will affect the response design, sampling design, and

analysis. The mmu obviously influences the partition of the ROI into map polygons because the mmudetermines the smallest area that is possible for a map polygon.

155 Pixel, block, and polygon assessment units are all arbitrary partitions of the ROI. Congalton and 156 Green (2008, p. 70) criticize using a single pixel as the spatial unit for an accuracy assessment stating that 157 a pixel is an arbitrary unit that does not have a meaningful relationship to most of the earth surface 158 features that are mapped. Although Congalton and Green (2008, p. 71) advocate using a block of pixels 159 as the spatial assessment unit, they acknowledge that a block of pixels is also arbitrary. The validity of 160 the results of an accuracy assessment does not depend on whether the assessment unit is a meaningful 161 entity, but instead depends on whether the area representation portrayed by the population error matrix is 162 meaningful. Regular size and shape spatial units (e.g. pixels) partition the landscape into convenient units 163 for analysis. Although surface characteristics such as land cover do not in reality typically conform to 164 such a spatial partition of the landscape, pixels still provide useful information regarding accuracy if the 165 partition formed by these units preserves the areas of agreement and disagreement of the population 166 defined by overlaying the reference map on the target map. For an area-based accuracy assessment, it is 167 the preservation of these areas of agreement and disagreement that is the critical trait required of the 168 spatial partition, and failure to preserve these areas is exacerbated by larger spatial units. Part of the 169 criticism of pixel-based assessments appears to be attributable to the failure to distinguish between area-170 based and per-object assessments. When conducting a per-object assessment in which a decision is made 171 whether each object is classified correctly, it would generally not be reasonable to examine a single pixel 172 to assess whether the majority of the area of a multiple-pixel object is correctly classified. But for an 173 area-based accuracy assessment in which accuracy is defined in terms of the location-specific area 174 representation of agreement and disagreement, a pixel assessment unit is a legitimate and practical option. 175

#### 176 2.1 Heterogeneity within spatial assessment units

For any choice of assessment unit, the decision of whether to treat these units as homogeneous or
heterogeneous (according to the reference and map classifications) will exert a strong influence on the

179 population. If the units are treated as homogeneous, a single class label can be assigned to each unit. In 180 contrast, if the units are regarded as heterogeneous, this heterogeneity must be accounted for in the class 181 labeling protocol. For example, the class information provided for a heterogeneous block may be a vector 182 of area proportions such as (0.8, 0.2) indicating that the area of the block is 80% class A and 20% class B. 183 Note that in this case the values 0.8 and 0.2 would still be regarded as constants in the design-based 184 framework, but the response is now a vector of observations instead of a single observation. Alternatively, 185 the class label information could be a fuzzy membership vector where membership is defined as the 186 degree to which the unit belongs to each of the classes (e.g., the vector (0.15, 0.85) indicates membership 187 of 0.15 in class A and membership of 0.85 in class B). When within-unit heterogeneity exists it is still an 188 option to impose an approach that essentially treats the assessment units as homogeneous. For example, a 189 unit with area proportions of 0.8 and 0.2 for classes A and B could be assigned a label of class A and the 190 unit subsequently treated as homogeneous in the analysis.

191 The definition of agreement and the data analysis will depend on this decision regarding 192 homogeneity of the assessment unit. If the assessment unit is treated as homogeneous in terms of both the 193 map and reference classification (i.e., the simplest case of a single map class and a single reference class) 194 then agreement exists if the map and reference labels match, and if the reference and map classes differ, it 195 is straightforward to specify the type of disagreement. A traditional error matrix analysis is readily 196 applied to the case of homogeneous assessment units. If the assessment unit is treated as heterogeneous, 197 analyses taking into account the mixed character of the unit are typically more complex and less familiar 198 to most practitioners than the error matrix analyses. Binaghi et al. (1999), Lewis and Brown (2001), 199 Pontius and Cheuk (2006) and Kuzera and Pontius (2008) are examples of analyses applicable to 200 heterogeneous assessment units in which the results are summarized by an error matrix and associated 201 accuracy measures. Another approach to quantify agreement between the target map and reference map 202 appropriate for heterogeneous assessment units is to estimate measures such as mean deviation, mean 203 absolute deviation, root mean square error, and correlation. Willmott (1982), Willmott and Matsuura 204 (2006), Ji and Gallo (2006), Pontius et al. (2008), and Riemann et al. (2010) critique these approaches and

suggest still other alternatives. It is not our intent to recommend particular measures or to review all analyses potentially applicable to heterogeneous assessment units. Instead, our purpose is to emphasize that an analysis different from the traditional error matrix approach will need to be applied in these cases.

209 2.2 Empirical demonstration of the effect of assessment unit on accuracy parameter values

210 To quantitatively illustrate the impact of the choices made for the spatial assessment unit and 211 response design, we constructed hypothetical populations for three regions. For the populations labeled 212 "Florida", the target map is the adjusted NLCD 1992 land cover of the United States (Fry et al. 2008) and 213 the reference map is the NLCD 2001 land cover (Homer et al. 2007). The other two regions, both located 214 in North Carolina (USA), are labeled "Fayetteville" and "Dare". For these two regions, the target map is 215 the NLCD 2001 land cover and the reference map is the NOAA C-CAP 2001 land cover. In reality, true 216 complete coverage reference maps do not exist for these regions, so available complete coverage maps 217 such as NLCD and C-CAP are used as hypothetical reference maps for the purpose of creating test 218 populations. The NLCD and C-CAP products have a 30-m x 30-m pixel resolution and each pixel has a 219 single map class label. The land-cover classes used in this analysis are Level I NLCD and C-CAP classes 220 water, urban, barren, forest, shrub, agriculture, and wetland. Detailed descriptions of the land-cover data 221 sources are available at www.mrlc.gov (NLCD) and www.csc.noaa.gov/digitalcoast/data/ccapregional/ 222 (C-CAP). The three target maps are shown in Figures S1-S3 of the supplemental online material. 223 Three assessment units are investigated in this analysis, a 30-m x 30-m pixel, a 3x3 pixel block, and a 224 map polygon, where the polygons are defined by the target map being assessed (i.e., a polygon is a 225 contiguous area of homogeneous mapped land cover, and contiguity is defined using the four neighboring 226 pixels). The pixel units are homogeneous in terms of both the map classification and reference 227 classification, but block units may be heterogeneous according to either the map or reference 228 classifications, and the map polygon assessment units may be heterogeneous according to the reference 229 classification. When heterogeneity within pixel blocks and polygons is present the mode class is used to 230 label the unit and agreement is defined based on comparing the mode map class and the mode reference

class. When more than one class qualifies as the mode for either the reference or map classification, a
label is assigned according to the following randomly ordered sequence: forest, wetland, shrub, urban,
barren, agriculture, and water. For example, if the map classification for a 3x3 block results in four pixels
of forest, four pixels of water, and one pixel of wetland, the class label assigned would be forest because
forest precedes water in the list used to decide the mode when ties occur. Other rules for resolving ties
could be constructed and would result in different populations. The rule we apply was chosen for
simplicity and so that similar blocks agree.

238 The population areas of disagreement and agreement obtained by overlaying the reference and target 239 maps provide visual evidence of the change in accuracy resulting from the choice of spatial unit (Figures 240 1-3). The differences in accuracy resulting from using different assessment units are summarized by 241 overall, user's, and producer's accuracies, with all accuracy parameter values (Table 2) computed from 242 the complete coverage map information. The three regions vary in how accuracy results change for 243 different spatial units. For Florida, accuracy does not change substantially for different spatial units, 244 whereas the Fayetteville region shows the largest differences in accuracy with the differences being 245 smaller between the pixel and block assessments relative to the differences between the pixel and polygon 246 assessments. The results for the Dare region are intermediate to the other two regions in that differences 247 in accuracy among the three spatial units exist but these differences are not as large as those observed for 248 the Fayetteville region. With the exception of some erratic differences observed for the very rare classes 249 (e.g., urban in the Dare population and water in the Fayetteville population) the largest difference in 250 accuracy between the pixel and block units occurs for user's accuracy of shrub in the Dare population 251 (11% difference). Changes in accuracy of 5% or more between the pixel and block results for user's 252 accuracy and producer's accuracy occur for several land-cover classes.

The accuracy results for the polygon assessment unit show greater variation from the pixel and block unit results. The differences are most dramatic for the Fayetteville population where overall accuracy is 12% higher for the polygon assessment compared to the pixel assessment. The class-specific accuracies of the polygon assessment can be much different from the pixel and block results. For example in the

Dare population, user's accuracy of forest is 85% for the polygon unit compared to 69% and 68% for the pixel and block units, and producer's accuracy of forest is 95% for the polygon unit compared to 79% and 77% for the pixel and block units. Producer's accuracy of shrub is another case in which class-specific accuracy is much higher for the polygon assessment than for the pixel or block assessment. For the Florida and Dare populations, the polygon accuracy results are generally slightly higher than the pixel and block results.

263 The results in Table 2 are intended to provide case study examples demonstrating the dependence of 264 accuracy results on the choice of spatial unit. Additional comparisons of accuracy parameter values based 265 on different spatial units for several additional populations are presented in the Appendix. We have not 266 attempted to discover or model how accuracy results for a particular ROI will vary depending on the 267 spatial assessment unit used. The size and shape of land-cover patches and accuracy of the classification 268 likely interact with other factors in a complicated manner to determine the sensitivity of the accuracy 269 results to the choice of assessment unit. A useful direction for future research would be to investigate if 270 the sensitivity of accuracy results to choice of spatial assessment unit can be modeled as a function of 271 landscape structure, classification accuracy, and other factors.

272

#### 273 2.3 Spatial Registration (Location) Error

274 Congalton and Green (2008) criticize the use of a pixel as an assessment unit, asserting that it is too 275 sensitive to positional errors that arise from fitting remotely sensed data to a map surface. Registration 276 errors affect all spatial units, and no spatial unit is entirely free of location error (see, for example, 277 McRoberts 2010). Conceptually, spatial misregistration is a "halo" around the assessment unit. For 278 example,  $a \pm 1$  pixel root mean square error (RMSE) means that the single pixel, the entire block of 279 pixels, or the entire polygon could be shifted  $\pm 1$  pixel along the x-dimension, y-dimension, or both 280 dimensions. Larger spatial units might be expected to be less sensitive to registration error because the 281 proportion of the total area of the spatial unit that is affected by location error decreases as the size of the 282 spatial unit increases. However, this decrease in sensitivity develops slowly as a function of increasing

283 size. For example, decreased sensitivity to spatial location error is modest for small pixel blocks as 284 illustrated by Figure 4. Panel A of Figure 4 depicts a 3x3 pixel block overlaid on a Landsat composite. 285 An interpreter collecting reference data would use the information in Panel A to locate the sample unit on 286 the reference medium. Once the sample unit had been located, the interpreter would use the reference 287 information (Panel C) to assign a reference label. If spatial location error exists, however, the sampling 288 unit could be shifted on the reference medium (Panel D). This misalignment of the map and reference 289 data is shown as a shift of the sample block location between Panel C and Panel D, with a change in the 290 mode reference class possibly resulting from the shift.

291 The ultimate impact of location error is determined by the change in the population and values of the 292 accuracy parameters resulting from location error relative to the population and parameter values when 293 location error is absent. To evaluate the effect of location error, we shifted all three reference maps by 294 one pixel down and one pixel to the left. The target map is unchanged by this shift in the reference map, 295 but the shift introduces location error and changes the population and values of the accuracy parameters 296 from the parameter values of the spatially aligned population (Table 2). The accuracy parameter values 297 resulting from the location-error impacted difference populations are shown in Table 3. In general, the 298 effect of location error is greatest for the pixel assessment unit followed by the block and polygon units. 299 For all three regions, location error produces a decrease in overall accuracy. The Florida population 300 showed the largest changes in accuracy attributable to location error with the decrease in overall accuracy 301 of 28% for the pixel unit, 20% for the block unit, and 12% for the polygon unit. Large decreases in class-302 specific accuracy were also observed for the Florida population. For example, user's accuracy for shrub 303 decreased by 38% for the pixel unit, 28% for the block unit, and 27% for the polygon unit, whereas 304 producer's accuracy for shrub decreased roughly 40%, 30%, and 20% for the pixel, block, and polygon 305 units. Relative to the Florida population, the Fayetteville and Dare populations generally had smaller 306 decreases in accuracy caused by shifting the reference map. More common land-cover classes generally 307 show smaller decreases in accuracy between the aligned and shifted reference maps. For example, urban 308 and agriculture are the two most common classes in the Fayetteville population and the decreases in

309 user's and producer's accuracies for these classes caused by location error are generally among the 310 smaller decreases observed (Table 3b). Similarly, forest and wetland are the two most common classes in 311 the Dare population and the decreases in accuracy for these classes are generally small, although the Dare 312 population is the least affected population in terms of results changing because of location error. 313 Use of an assessment unit larger than a pixel to mitigate the impact of location error is supported to 314 some degree by these results (Table 3). However, the decreases in overall accuracy for both the block and 315 polygon units provides evidence that these units are also susceptible to substantial errors in accuracy 316 attributable to spatial misregistration between the map and reference locations. The effect of location

317 error when using blocks or pixels is also evident for user's and producer's accuracies since most land-

318 cover classes have lower values for these parameters for the spatially misaligned population.

319

#### 320 **3.** The impact of spatial assessment unit on sampling design and estimation

#### 321 *3.1 Defining a sampling universe and frame*

322 Constructing a sampling design requires specifying a sampling universe, defined as the set of spatial 323 units that form a partition of the ROI. In practice, a universe of pixels, blocks of pixels, or map polygons 324 is straightforward to construct. Because the universe must be spatially exhaustive, it is not valid to 325 exclude heterogeneous areas as is sometimes done to avoid location error issues. The pixels, blocks, or 326 polygons forming the universe can be represented by a "list frame", defined as a list of all such units in 327 the ROI along with a spatial address for each unit (e.g., spatial coordinates or an identification number 328 unique to each unit). Because pixels and blocks provide a partition of the ROI independent of the map or 329 reference classification, it is possible to construct a list frame of pixels or blocks before the map is finalized. A list frame of map polygons could be readily produced from either a final or preliminary map, 330 331 although using a preliminary map is a questionable practice because not all of these map polygons will 332 exist in the final map. If polygons are defined by the reference classification, it is not feasible to construct 333 a complete list frame of the universe of reference polygons because a census of the reference 334 classification would be required. Protocols for implementing basic sampling designs such as simple

335 random, stratified random, systematic, and cluster sampling from a list frame are described in Cochran 336 (1977), Lohr (1999), and Särndal et al. (1992). Stehman (1999, 2009) provides a general overview of 337 sampling designs applicable to accuracy assessment.

- 338
- 339

#### *3.2 Stratified sampling and cluster sampling*

340 Two important considerations when choosing the sampling design are whether to group the 341 spatial assessment units into strata to control the sample size allocated per stratum (for the purpose of 342 decreasing standard errors of class-specific accuracy estimates) and whether to group the units into 343 clusters to spatially constrain the sample within these clusters (for the purpose of decreasing costs 344 associated with travel time to field sites or costs associated with the number of aerial photographs or very 345 high resolution images required to obtain the reference data). The choice of a pixel, block, or polygon 346 assessment unit has ramifications on how such stratified and cluster sampling designs would be 347 implemented. We will address the stratification and clustering considerations separately. 348 Typically in accuracy assessment sampling designs the strata correspond to the mapped area of each 349 class (e.g., all area mapped as forest is one stratum) and the sample sizes allocated to rare class strata are 350 chosen to achieve specified standard errors of the user's accuracy estimates. If each assessment unit has a 351 single map class, as is the case for a pixel or a map polygon, it is a simple matter to assign each unit to a 352 single stratum corresponding to the unit's map class. Stratification is more complicated for a block 353 assessment unit because not all blocks are comprised of a single map class, and a protocol stating how to 354 assign each heterogeneous block to a stratum must be specified. A variety of assignment rules could be 355 envisioned to create the stratification of blocks (e.g., a block could be assigned to a stratum based on the 356 most common map class within the block or based on the class associated with the center of the block),

357 but the effectiveness of different stratum assignment rules has not been investigated.

358 A cluster is a group of pixels, blocks, or polygons that is treated as a single entity in the sampling 359 protocol. Pixel and block assessment units can be easily grouped into regularly-shaped clusters so pixel 360 and block units conveniently fit the nested structure of cluster sampling. Forming clusters of polygons is

361 more cumbersome than forming clusters of pixels or clusters of blocks. Polygons vary in size and shape, 362 so there is no intuitive way to group polygons to form clusters so that the clusters are approximately 363 uniform in size and shape. Polygons lack the regular size and adjacency features possessed by pixels and 364 blocks that allow easy nesting of pixels or blocks within same size clusters. The advantage of cluster 365 sampling to constrain the sample is diminished for polygons because of the variation in the size of the 366 clusters formed. For example, a cluster of four very large polygons covers a much larger area than a 367 cluster of four small polygons. Analysis of cluster sampling designs is considerably simpler when the 368 clusters are equal in size so analyzing a cluster sample of polygons will be more complex relative to 369 analyzing a cluster sample of pixel or block assessment units.

Note that a cluster of pixels appears at first glance to be the same as a block assessment unit but a cluster of pixels is treated differently from a block unit. The response design associated with a block assessment unit does not produce a reference label for each pixel within the block but rather assigns the reference classification to the block as a single entity. In contrast, a cluster comprised of pixels would include a reference classification for each pixel in the cluster, so a pixel remains the assessment unit within the cluster, and the cluster is the primary sampling unit (Stehman 1997).

376

#### 377 3.3 Point sampling to select pixel, block, or polygon units

378 An alternative to selecting the sample from a list frame is spatial point sampling. In this protocol, a 379 spatial sample of points is first selected within the ROI and the pixel, block, or polygon within which a 380 sample point falls is selected into the sample. Point sampling is effectively the same as sampling from a 381 list frame if the spatial units partitioning the ROI are all equal in area and have the same shape, as would 382 be the case for pixel or block assessment units, and the decision of whether to use a list frame or point 383 sampling approach with pixel or block assessment units would be made on the basis of convenience of 384 implementation. In contrast, for a polygon assessment unit, the point sampling protocol will select 385 polygons into the sample with probability proportional to polygon area, so larger polygons will have a

higher probability of being selected (Figure 5). These unequal inclusion probabilities must be accountedfor in the analysis.

The point sampling approach provides a way to select a sample when it is impractical to construct a list frame. For example, if the spatial unit is a polygon defined by the reference classification (i.e., a "reference polygon"), a list frame of all such polygons is not available. To obtain a sample of reference polygons by point sampling, the sample point locations are selected, and the reference polygon that contains each sample point is delineated based on the reference classification (Figure 5). Only those reference polygons intercepted by sample points would need to be identified, so it is no longer necessary to create a list of all reference polygons comprising the population.

395 Stratified and cluster sampling would not be viable options when point sampling is used to select 396 reference polygons. Stratification requires assigning all reference polygons in the entire ROI to strata 397 based on each polygon's reference class label, and obtaining this information would be impractical 398 because it is tantamount to a census of reference polygons. Two-phase sampling (Cochran 1977) in which 399 the stratification assignment is required only for a large first-phase sample instead of a census ameliorates 400 this practical disadvantage associated with a reference polygon assessment unit. Cluster sampling of 401 reference polygons is not viable because all reference polygons within a cluster would need to be 402 delineated, and this places an impractical burden on reference data collection.

403

404 *3.4 Estimation* 

The values of the population parameters resulting from the choice of response design and assessment unit are estimated from the sample. The sampling design has no effect on these parameter values because in the design-based inference framework the value of a population parameter remains fixed regardless of the particular sample selected. It is the <u>estimate</u> of that fixed parameter value that changes depending on the sample and sampling design. It is critical to recognize that the value of a target parameter (e.g., overall, user's, or producer's accuracy) is determined by the choice of spatial assessment unit and response design (see Section 2), not by the sampling design. The specific formula used to

412 estimate a parameter will differ depending on the sampling design and the precision of an estimate will413 vary for different sampling designs.

414 The theory of probability sampling ensures the availability of an estimator that is either unbiased 415 or consistent for the parameter of interest. If the accuracy estimator is a Horvitz-Thompson estimator 416 (Stehman 2001, p. 728), it is an unbiased estimator (Särndal et al. 1992, p. 43). Sometimes an unbiased 417 estimator is not available, as may occur when estimating a ratio of two parameters (e.g., producer's 418 accuracy is the estimated area of agreement for a specific class divided by the estimated area of that class 419 according to the reference classification). In such cases, a consistent estimator can be constructed 420 (Särndal et al. 1992, Sec. 5.3; Overton and Stehman 1995), where a consistent estimator is one in which 421 "... the sampling distribution of the estimator can be considered tightly concentrated around  $\tau$  [the 422 parameter], when n [the sample size] is large enough" (Särndal et al. 1992, p. 166). Thus an unbiased 423 estimator guarantees that the parameter of interest is estimated correctly on average (averaging over all 424 possible samples that could be selected), and a consistent estimator (although not necessarily unbiased) 425 ensures that the estimate for a particular sample will not stray too far from the true value of the parameter. 426 The values of the parameters targeted by the sample-based estimators depend on the population 427 created by the overlay of the target map and reference map and the definition of agreement specified by 428 the response design (see Section 2). Although sampling theory exists to support unbiased or consistent 429 estimation of any population parameter, the specific estimator formulas for some combinations of 430 parameter and sampling design may need to be derived. Typically estimators used in accuracy 431 assessment are presented only for simple random sampling (Congalton and Green 2008) and different 432 estimators are needed for other sampling designs (e.g., Card 1982, Stehman 1996, Stehman and Foody 433 2009). Deriving variance estimators may also be necessary, again most likely when the design is not 434 simple random sampling.

To reiterate the key message of this section, once the population and values of the parameters are determined by the choice of response design and spatial unit, the sampling design and data analysis protocol can provide statistically defensible estimates of these parameters. Any probability sampling

design will permit either an unbiased or a consistent estimator of a parameter, so there is no distinction
among sampling designs relevant to this feature of estimation. Different sampling designs will yield
different values for the variance of an estimator of a particular accuracy parameter, so the comparison of
sampling designs should be on the basis of the variance of the estimator and not bias (or consistency) of
the estimator.

443

#### 444 **4. Additional Considerations**

445 *4.1 Block assessment units* 

446 Choosing a 3x3 pixel block as the assessment unit and defining agreement based on the mode class of 447 the block results in the map population assessed being different from the map provided to users. Defining 448 agreement based on the mode class of a 3x3 pixel block implies assessment of a map that has been 449 "filtered" to produce a classification at the support of a block (e.g., an area of 900 m<sup>2</sup> for a 30-m x 30-m 450 pixel). For example, suppose a 3x3 pixel block has five of the nine pixels labeled as forest according to 451 the reference classification, and five of the nine labeled as forest according to the map, but only one pixel 452 in common is forest according to both the map and reference classification. In terms of total forest area of 453 the block, agreement is 55% (5 out of 9), yet the overlay of the target and reference maps for the block 454 would show only 11% agreement of common overlapping area of forest (1 out of 9 pixels). The mode 455 class for both the map and reference classification is forest, so if agreement is defined by comparing the 456 mode class, the block would be classified correctly even though only a single pixel has forest for both the 457 map and reference classifications.

458 Czaplewski (2003) provides a pointed criticism of accuracy assessments using data aggregated to a 459 block level when the map provided to users is not similarly aggregated. The mismatch between the 460 spatial support at which the accuracy assessment is conducted and the spatial support of the map brings 461 into question the utility of the accuracy results obtained from such an approach. In contrast, response 462 designs that use a single pixel as the spatial unit assess the map exactly as it is distributed to users. These

463 considerations reinforce the importance of clearly specifying the population that is the targeted objective464 of the accuracy assessment.

465

#### 466 4.2 Polygon assessment units

467 Use of a map polygon assessment unit introduces concerns not present for pixel or block assessment 468 units because the utility of a map polygon for accuracy assessment is inseparably dependent on the target 469 map. A map polygon may become obsolete for accuracy assessment if the map undergoes a revision that 470 changes the map polygons forming the original partition of the ROI that existed when the sample of map 471 polygons was selected. The most intuitive example of revision would be aggregation of classes to form a 472 simplified legend (e.g., Anderson et al. (1976) Level II to Level I) and a re-analysis of the data based on 473 the aggregated classification (e.g., Pontius and Malizia 2004). Class aggregation re-draws the map 474 polygons, so a sample of map polygons based on a Level II classification would not necessarily provide a 475 sample of map polygons that existed for the Level I classification. Because the sampling and response 476 designs are directly linked to the specific map polygons partitioning the ROI, it would generally not be 477 possible to use a partition based on Anderson Level II polygons to estimate accuracy for the map of Level 478 I polygons. Pixels and blocks of pixels retain their utility even if the map classes are aggregated or the 479 map is otherwise revised because their spatial boundaries remain defined and fixed despite label changes. 480 Because polygon assessment units vary in size, stratifying the population by polygon size merits 481 consideration. For example, for the Florida population, the map polygons (obtained from the NLCD 482 1992) range in size from 0.09 ha (a single pixel) to 139.95 ha, and the largest 2.3% of the polygons 483 account for 50% of the total area. For the reference map (obtained from the NLCD 2001), the range in 484 polygon size is 0.09 ha to 231.93 ha, with the largest 2.0% of the polygons accounting for 50% of the 485 total area. An equal probability sampling design (e.g., simple random, systematic, or stratified random 486 with proportional allocation) will result in a sample with the same size distribution of polygons as the 487 population, so a large proportion of the sample will consist of small polygons. Stratification by polygon 488 size could be used to increase the proportional representation of larger polygons in the sample, but it is

489 unclear whether there is an advantage to increasing the sample size of large polygons as no studies have 490 examined the impact of different sampling design choices on the standard errors of accuracy estimates 491 obtained from a polygon-based assessment. Further, stratification by polygon size would likely be 492 combined with stratification by map class and this two-way stratification increases the overall complexity 493 of the sampling design and analysis.

494 Very large polygons also introduce challenging problems for the response design. Obtaining the 495 reference classification for the entire area of a very large polygon may be too expensive or impractical, so 496 a portion of the polygon may be sampled and the reference characteristics of the polygon estimated from 497 the sample. There is little research on how these factors affect accuracy assessment results.

498

514

#### 499 **5.** Summary

500 The three most commonly used spatial units for accuracy assessment are a pixel, a block of pixels, 501 and a map polygon. While a universally best spatial assessment unit does not exist, the choice of spatial 502 unit has broad implications on the conduct and outcome of the assessment. The results of a map accuracy 503 assessment depend on the spatial unit chosen to serve as the basis of the assessment because different 504 spatial units lead to different populations and values of the accuracy parameters. The population 505 perspective (see Section 2) provides a rigorous conceptual framework for evaluating the impact of the 506 choice of spatial unit. Specifically, the population of interest in accuracy assessment may be viewed as 507 the result of overlaying a complete coverage reference map with the map to be evaluated and quantifying 508 the class-specific areas of agreement and disagreement between the reference map and the target map. 509 The values of the accuracy parameters computed from this population are the quantities estimated from 510 the sample of reference data. Greater awareness of this population framework will clarify the 511 ramifications of the choice of spatial unit on the outcome of an accuracy assessment (Table 1). 512 The focus of this article is area-based accuracy assessments in which the area of each land-cover type 513 correctly classified and the area incorrectly classified by type of misclassification are the primary data for

20

describing accuracy. These areas are summarized by an error matrix and the accuracy parameters derived

from the error matrix. The difference map created by overlaying a complete coverage reference classification on the target map generates the population (i.e., the per-class areas of agreement and the areas of each class-specific type of disagreement). The ROI is then partitioned by the spatial unit chosen for the accuracy assessment.

519 As the smallest spatial assessment unit, a pixel best preserves the population areas of agreement and 520 disagreement when the ROI is partitioned. Many accuracy assessments have been conducted using a 521 pixel as the assessment unit. Therefore, specific details of the sampling design, response design, and 522 analysis protocols associated with implementing a pixel-based assessment have been extensively 523 developed and applied. A pixel-based assessment easily accommodates sampling designs employing 524 strata or clusters, whereas blocks are less amenable to stratification and polygons are less practical to use 525 in cluster sampling. The traditional error matrix analyses are readily implemented for a pixel-based 526 assessment.

527 Blocks and polygons are less likely than pixels to be homogeneous so the response design and 528 analysis protocols must be more complex to account for within-unit heterogeneity. A common practice 529 when using a block or polygon assessment unit is to revert to the protocols developed for a pixel-based 530 assessment and to assume (although rarely explicitly stated) that the block or polygon units are 531 homogeneous. In the likely case that within-unit heterogeneity is present, the mode class is often 532 assigned as the class label and this labeling protocol changes the areas of agreement and disagreement 533 defining the population and correspondingly changes the values of the accuracy parameters. This is a 534 critical feature of accuracy assessment that must be recognized. Assessments based on block or polygon 535 units generally do not preserve the areas of agreement and disagreement that would be obtained by 536 overlaying the unpartitioned target map and the unpartitioned reference map.

The accuracy results for the example populations (Table 2) illustrated a range of outcomes from little change in accuracy with different spatial units (Florida population) to substantial differences in accuracy with different spatial units (Fayetteville and Dare populations). The impact of choice of spatial unit on the results of an accuracy assessment is not only important for single map assessments, but also for

541 comparative studies of accuracy. For example, a researcher evaluating the performance of a new classifier 542 by comparison against an existing classifier should consider the spatial unit used in the assessment of the 543 two classifiers (and perhaps use the same unit to avoid confounding sources of uncertainty).

The choice of spatial unit to serve as the basis of an accuracy assessment is a critical decision. The spatial unit must be chosen with the understanding that the population and therefore the values of the accuracy parameters describing the population are determined by the spatial unit in combination with the response design. Further, the sampling design must be appropriate for the spatial unit chosen. Sampling designs using strata and clusters that are commonly and easily implemented for a pixel unit are more cumbersome to implement when using block or polygon units. Better recognition of the impacts of the choice of spatial unit and the advantages and disadvantages of each unit will lead to better accuracy

assessment methodology and improve the validity of the results.

552

#### 553 Acknowledgments

554 We thank four anonymous reviewers for their helpful and constructive comments. SVS acknowledges

555 funding support through an Intergovernmental Personnel Agreement (IPA) from the U.S. Geological

556 Survey Earth Resources Observation and Science (EROS) Center, Sioux Falls, South Dakota. This article

- has not been subject to review by the U.S. Environmental Protection Agency or U.S. Geological Survey
- and does not necessarily reflect the views of either agency.

559

#### 560 **References**

- 561
- Anderson, J. R., Hardy, E. E., Roach, J. T., & Witmer, R. E. (1976). A land use and land cover
  classification system for use with remote sensor data, U.S. Geological Survey Professional Paper 964,
  Washington, DC: U.S. Geological Survey, 28 p.
- 566 Binaghi, E., Brivio, P. A., Ghezzi, P., & Rampini, A. (1999). A fuzzy set-based accuracy assessment of 567 soft classification, *Pattern Recognition Letters*, 20, 935-948.
- 568

565

Card, D. H. (1982). Using known map category marginal frequencies to improve estimates of thematic
 map accuracy. *Photogrammetric Engineering and Remote Sensing*, 48, 431-439.

572 Cochran, W. G. (1977). Sampling Techniques (3rd edn.). New York: John Wiley & Sons.

573

- 574 Congalton, R. G., & Green, K. (1999). *Assessing the Accuracy of Remotely Sensed Data: Principles and* 575 *Practices.* Boca Raton, FL: CRC Press.
- 576
- 577 Congalton, R. G., & Green, K. (2008). Assessing the Accuracy of Remotely Sensed Data: Principles and 578 Practices (2<sup>nd</sup> edn.). Boca Raton, FL: CRC Press.
- 579
- 580 Czaplewski, R. L. (2003). Accuracy assessment of maps of forest condition: Statistical design and
- 581 methodological considerations. In *Remote Sensing of Forest Environments: Concepts and Case Studies*
- 582 (Wulder, M. A., & Franklin, S. E., eds.). Boston: Kluuwer Academic Publishers, pp. 115-140. 583
- Dungan, J. L. (2006). Focusing on feature-based differences in map comparison, *Journal of Geographical Systems*, 8, 131-143.
- 586
- 587 Fry, J. A., Coan, M. J., Homer, C. G., Meyer, D. K., & Wickham, J. D. (2009). Completion of the
- 588 National Land Cover Database (NLCD) 1992-2001 Land Cover Change Retrofit Product. U. S.
- Geological Survey Open-File Report 2088-1379 (available from <u>http://pubs.usgs.gov/of/2008/1379</u>,
   accessed April 2011).
- 591
- Hargrove, W. W., Hoffman, F. M., & Hessburg, P. F. (2006). Mapcurves: a quantitative method for
   comparing categorical maps. *Journal of Geographical Systems*, 8, 187-208.
- 594

Homer, C., Dewitz, J., Fry, J., Coan, M., Hossain, N., Larson, C. et al. (2007). Completion of the 2001
National Land Cover Database for the conterminous United States. *Photogrammetric Engineering and Remote Sensing*, *73*, 337-341.

- 597 598
- Janssen, L. L. F., & van der Wel, F. J. M. (1994). Accuracy assessment of satellite derived land-cover
  data: A review. *Photogrammetric Engineering and Remote Sensing*, 60, 419-426.
- Ji, L., & Gallo, K. (2006). An agreement coefficient for image comparison, *Photogrammetric Engineering and Remote Sensing*, 72:823-833.
- 604

Kuzera, K., & Pontius, R. G., Jr. (2004). Categorical coefficients for assessing soft-classified maps at
 multiple resolutions. *In* Proceedings of the joint meeting of the 15th annual conference of The
 International Environmetrics Society and the 6th annual symposium on Spatial Accuracy Assessment in

- 608 the Natural Resources and Environmental Sciences, 28 June 1 July 2004, Portland, ME. 609
- Lewis, H. G., & Brown, M. (2001). A generalized confusion matrix for assessing area estimates from
  remotely sensed data. *International Journal of Remote Sensing*, 22, 3223-3235.
- 613 Lohr, S. L. (1999). Sampling: Design and Analysis. New York: Duxbury Press.
- 615 McRoberts, R. E. (2010). The effects of rectification and Global Positioning System errors on satellite 616 image-based estimates of forest area, *Remote Sensing of Environment*, *114*, 1710–1717.
- 617
  618 Overton, W. S., & Stehman, S. V. (1995). The Horvitz-Thompson theorem as a unifying perspective for
  619 probability sampling: with examples from natural resource sampling. *American Statistician*, 49, 261-268.

620 621 Pontius, R. G., Jr., & Cheuk, M. L. (2006). A generalized cross-tabulation matrix to compare soft-

- 622 classified maps at multiple spatial resolutions. *International Journal of Geographical Information*
- 623 Science, 20, 1-30.
- 624

- 625 Pontius, R. G., Jr., & Malizia, N. R. (2004). Effect of category aggregation on map comparison. In
- Lecture Notes in Computer Science 3234: Proceedings of the Third International Conference, GIScience
  2004 (Egenhofer, M. J., Freksa, C., & Miller, H. J., eds). Adelphi, Maryland (USA): Springer, pp. 251268.
- 628 629
- 630 Pontius, R. G., Jr., Thontteh, O., & Chen, H. (2008). Components of information for multiple resolution
- 631 comparison between maps that share a real variable. *Environmental and Ecological Statistics*, *15*, 111 632 142.
- 633
- 634 Power, C., Simms, A., & White, R. (2001). Hierarchical fuzzy pattern matching for the regional
- 635 comparison of land use maps. *International Journal of Geographical Information Science*, *15*, 77-100.636
- Richards, J. A. (1996). Classifier performance and map accuracy, *Remote Sensing of Environment*, 57, 161-166.
- 639
- 640 Riemann, R., Wilson, B. T., Lister, A., & Parks, S. (2010). An effective assessment protocol for
- 641 continuous geospatial datasets of forest characteristics using USFS Forest Inventory and Analysis (FIA)
  642 data. *Remote Sensing of Environment*, 114, 2337-2352.
- 643
- 644 Särndal, C. E., Swensson, B., & Wretman, J. (1992). *Model-Assisted Survey Sampling*. New York:
  645 Springer-Verlag.
- 646
- Smith, S. V., Renwick, W. H., Bartly, J. D., & Buddemeier, R. W. (2002). Distribution and significance
  of small, artificial water bodies across the United States landscape. *The Science of the Total Environment*, 299, 21-36.
- Stehman, S. V. (1996). Estimating the Kappa coefficient and its variance under stratified random
  sampling. *Photogrammetric Engineering and Remote Sensing*, *62*, 401-407.
- 654 Stehman, S. V. (1997). Estimating standard errors of accuracy assessment statistics under cluster 655 sampling. *Remote Sensing of Environment*, 60, 258-269.
- 657 Stehman, S. V. (1999). Basic probability sampling designs for thematic map accuracy assessment,
  658 *International Journal of Remote Sensing*, 20, 2423-2441.
- 660 Stehman, S. V. (2001). Statistical rigor and practical utility in thematic map accuracy assessment, 661 *Photogrammetric Engineering and Remote Sensing*, 67, 727-734.
- 662

- Stehman, S. V. (2009). Sampling designs for accuracy assessment of land cover. *International Journal of Remote Sensing*, 30, 5243-5272.
- Stehman, S. V., & Czaplewski, R. L. (1998). Design and analysis for thematic map accuracy assessment:
  Fundamental principles. *Remote Sensing of Environment*, 64, 331-344.
- 668 669 Stehman, S.V., & Foody, G. M. (2009). "Accuracy Assessment" (Chapter 21). In *The SAGE Handbook of* 
  - *Remote Sensing* (Warner, T. A., Nellis, M. D., & Foody, G. M. eds.). London: Sage Publications Ltd.,
     pp. 297-309.
  - 672
  - 673 Story, M., & Congalton, R.G. (1986). Accuracy assessment: a user's perspective. *Photogrammetric*
  - 674 Engineering and Remote Sensing, 52, 397-399.
  - 675

- 676 Strahler, A. H., Boschetti, L., Foody, G. M., Friedl, M. A., Hansen, M. C., Herold, M., et al. (2006).
- 677 Global land cover validation: Recommendations for evaluation and accuracy assessment of global land
- 678 cover maps, EUR 22156 EN - DG, Office for Official Publications of the European Communities, Luxembourg, 48 pp.
- 679 680
- 681 White, R. (2006). Pattern based map comparisons. Journal of Geographical Systems, 8, 145-164.
- 682 683 Willmott, C. (1982). Some comments on the evaluation of model performance. Bulletin of the American
- 684 Meteorological Society, 63, 1309-1313.
- 685
- 686 Willmott, C. J., & Matsuura, K. (2006). On the use of dimensioned measures of error to evaluate the
- 687 performance of spatial interpolators. International Journal of Geographical Information Science, 20, 89-688 102.
- 689

690 **Table 1.** Useful concepts and results to guide selection of an accuracy assessment spatial unit. 691 692 **General**: The three primary components of an accuracy assessment are the response design, sampling 693 design, and analysis (Stehman and Czaplewski 1998). The choice of spatial assessment unit (e.g., pixel, 694 block of pixels, or polygon) must be considered in terms of the ramifications on all three components. 695 696 **Defining the population and accuracy parameters** 697 698 The population that determines the values of the accuracy parameters targeted by the 0 699 assessment is conceptualized as resulting from overlaying a complete coverage reference 700 classification (i.e., a reference map) and the target map to be evaluated. The population 701 may be viewed as a difference map resulting from this overlay showing where the map 702 and reference classifications agree and where they disagree. The corresponding areas of 703 class-specific agreement and class-specific disagreement may be obtained from the 704 difference map. The difference map will depend on the partition of the region of interest 705 (ROI) created from the spatial unit chosen for the assessment. Different populations will 706 yield different values of the accuracy parameters (e.g., different overall accuracy values). 707 708 For an area-based accuracy assessment, the per-class area of agreement and area of 0 709 disagreement are summarized by a population error matrix where the cells of the error 710 matrix represent the proportion of area of agreement for each class and the proportion of 711 area misclassified for each type of error. The population and accuracy parameter values 712 obtained from the population error matrix will differ depending on the choice of spatial 713 unit. Changing the spatial unit changes the population and consequently the results of the 714 accuracy assessment. 715 716 • Pixels, blocks of pixels and polygons are all arbitrary spatial units. The validity of an 717 accuracy assessment does not depend on whether the spatial assessment unit is a real 718 surface feature of the earth but instead depends on the area representation of agreement 719 and disagreement resulting from use of the spatial unit to partition the ROI. 720 721 The population and values of the accuracy parameters are not affected by the choice of 0 722 sampling design. 723

724		
725	• Response	design – determining the reference classification and definition of agreement
726		
727	0	Heterogeneity within the spatial unit requires specification of how heterogeneity will be
728		accommodated in the labeling protocol and definition of agreement. These decisions will
729		have a substantial impact on the population and values of the accuracy parameters. As a
730		general rule, the likelihood of heterogeneity within a spatial unit will increase with size of
731		the unit.
732		
733	0	Map polygons do not necessarily represent real earth surface features and use of map
734		polygons does not support reporting of accuracy results for different levels of class
735		aggregations. The potentially large variation in polygon size creates practical challenges
736		to the response design protocol and also motivates consideration of stratifying the
737		sampling design by polygon size. Little research has been conducted exploring how
738		variation in polygon size affects reference data collection, sampling, and analysis.
739		
740	0	Polygons defined by the reference classification ("reference polygons") are appealing
741		because they represent real objects of interest but present practical challenges when
742		constructing the response and sampling designs. Methodological developments are
743		needed before reference polygons can be considered a viable option.
744		
745	• Sampling	design and estimation
746		
747	0	The sampling design requires specifying the universe of all spatial units forming a
748		partition of the ROI. The spatial units making up the universe must be non-overlapping
749		and spatially exhaustive.
750		
751	0	For a stratified design, each spatial unit must be assigned to one and only one stratum.
752		Block assessment units may be internally heterogeneous and this complicates the
753		protocol for assigning each block to a stratum.
754		
755	0	Polygons are less amenable than pixel or block assessment units for the purpose of cluster
756		sampling because polygons have no natural grouping or nesting structure to form
757		clusters.

758		
759	0	It is possible to construct an unbiased or a consistent estimator for any accuracy
760		parameter of interest when a probability sampling design is implemented and working
761		within the design-based inference framework. Consequently, it is pointless to consider
762		comparing sampling designs on the basis of bias or consistency of estimators derived for
763		different designs.
764		
765	0	Different sampling designs result in different precision (i.e., variance) for the estimator of
766		a given accuracy parameter, so it is meaningful to compare sampling designs on the basis
767		of variance of the accuracy estimators.
768		
769	Location	or registration error (map and reference data are not spatially aligned)
770		
771	0	Registration errors affect all spatial units. Spatial misregistration can be conceptualized
772		as a "halo" around the assessment unit.
772 773		as a "halo" around the assessment unit.
772 773 774	0	as a "halo" around the assessment unit. The ultimate impact of location error is observed in the change in the population and
<ul><li>772</li><li>773</li><li>774</li><li>775</li></ul>	0	as a "halo" around the assessment unit. The ultimate impact of location error is observed in the change in the population and values of the accuracy parameters relative to the parameter values of a population
<ul> <li>772</li> <li>773</li> <li>774</li> <li>775</li> <li>776</li> </ul>	0	as a "halo" around the assessment unit. The ultimate impact of location error is observed in the change in the population and values of the accuracy parameters relative to the parameter values of a population determined by overlaying perfectly spatially aligned reference and target maps (i.e.,
<ul> <li>772</li> <li>773</li> <li>774</li> <li>775</li> <li>776</li> <li>777</li> </ul>	0	as a "halo" around the assessment unit. The ultimate impact of location error is observed in the change in the population and values of the accuracy parameters relative to the parameter values of a population determined by overlaying perfectly spatially aligned reference and target maps (i.e., location error is absent).
<ul> <li>772</li> <li>773</li> <li>774</li> <li>775</li> <li>776</li> <li>777</li> <li>778</li> </ul>	O	as a "halo" around the assessment unit. The ultimate impact of location error is observed in the change in the population and values of the accuracy parameters relative to the parameter values of a population determined by overlaying perfectly spatially aligned reference and target maps (i.e., location error is absent).
<ul> <li>772</li> <li>773</li> <li>774</li> <li>775</li> <li>776</li> <li>777</li> <li>778</li> <li>779</li> </ul>	0	as a "halo" around the assessment unit. The ultimate impact of location error is observed in the change in the population and values of the accuracy parameters relative to the parameter values of a population determined by overlaying perfectly spatially aligned reference and target maps (i.e., location error is absent). Pixel-based assessments are generally more sensitive to location error than are block- and
<ul> <li>772</li> <li>773</li> <li>774</li> <li>775</li> <li>776</li> <li>777</li> <li>778</li> <li>779</li> <li>780</li> </ul>	0	<ul> <li>as a "halo" around the assessment unit.</li> <li>The ultimate impact of location error is observed in the change in the population and values of the accuracy parameters relative to the parameter values of a population determined by overlaying perfectly spatially aligned reference and target maps (i.e., location error is absent).</li> <li>Pixel-based assessments are generally more sensitive to location error than are block- and polygon-based assessments as evidenced by larger changes in the values of the accuracy</li> </ul>
<ul> <li>772</li> <li>773</li> <li>774</li> <li>775</li> <li>776</li> <li>777</li> <li>778</li> <li>779</li> <li>780</li> <li>781</li> </ul>	0	as a "halo" around the assessment unit. The ultimate impact of location error is observed in the change in the population and values of the accuracy parameters relative to the parameter values of a population determined by overlaying perfectly spatially aligned reference and target maps (i.e., location error is absent). Pixel-based assessments are generally more sensitive to location error than are block- and polygon-based assessments as evidenced by larger changes in the values of the accuracy parameters between spatially aligned and spatially unaligned reference and target maps.
<ul> <li>772</li> <li>773</li> <li>774</li> <li>775</li> <li>776</li> <li>777</li> <li>778</li> <li>779</li> <li>780</li> <li>781</li> <li>782</li> </ul>	0	as a "halo" around the assessment unit. The ultimate impact of location error is observed in the change in the population and values of the accuracy parameters relative to the parameter values of a population determined by overlaying perfectly spatially aligned reference and target maps (i.e., location error is absent). Pixel-based assessments are generally more sensitive to location error than are block- and polygon-based assessments as evidenced by larger changes in the values of the accuracy parameters between spatially aligned and spatially unaligned reference and target maps. However, location error can still have a considerable effect on accuracy results when

Table 2. Parameter values for user's and producer's accuracies resulting from different spatial assessment units, pixel, 3x3 pixel block, and map polygon. If a block or polygon is heterogeneous in terms of the map or reference classification, the mode class is used and agreement is defined as a match between the mode map class and the mode reference class. Parameter values for block and polygon units are reported as deviations from the pixel-based parameter value. Positive deviations indicate higher accuracy for the block- or polygon-based assessment, and negative deviations indicate lower accuracy for the block- or polygon-based assessment. Classes are ordered by decreasing percent of area mapped.

- 791
- a) Florida

793			Mean						
794		Мар	Patch	User's	Accura	<u>cy (%)</u>	Produc	er's Acc	uracy (%)
795	Class	Area (%)	Size(ha)	Pixel	Block	Polygon	Pixel	Block	Polygon
796	Wetland	20	16.6	100	0	0	55	-6	-6
797	Agriculture	20	4.8	99	0	0	98	0	+1
798	Forest	19	3.8	60	-5	-3	86	-2	-1
799	Urban	18	20.1	100	0	0	95	-2	+5
800	Shrub	16	2.3	85	-3	0	92	-1	+5
801	Water	7	12.4	92	0	+6	100	0	0
802	Barren	<1	0.9	100	0	0	100	0	0
803	Overall accu	uracy		86	-3	-1			
804									
805	b) Fayettevi	lle							
806			Mean						
807		Map	Patch	<u>User'</u>	s Accura	<u>acy (%)</u>	Produc	er's Acc	uracy (%)
808	Class	Area (%)	Size(ha)	Pixel	Block	Polygon	Pixel	Block	Polygon
809	Agriculture	32	6.8	86	+2	+5	83	+2	+6
810	Urban	31	29.2	72	+4	+18	85	+2	+5
811	Forest	17	2.6	69	-1	+16	79	-2	+16
812	Wetland	10	3.3	85	0	+12	65	0	+14
813	Shrub	9	0.9	68	-7	+6	54	-5	+29
814	Water	<1	1.7	89	+2	+5	61	+16	+20
815	Overall accu	uracy		77	+1	+12			
816									

818	c) Dare									
819			Mean							
820		Map	Patch	User	User's Accuracy (%)			Producer's Accuracy (%)		
821	Class	Area (%)	Size(ha)	Pixel	Block	Polygon		Pixel	Block	Polygon
822	Forest	45	24.1	90	+1	+8		91	-1	-1
823	Wetland	34	20.1	87	0	0		87	+1	+8
824	Shrub	9	2.3	71	+11	+11		77	0	+1
825	Agriculture	9	5.3	77	+6	+7		81	+1	-4
826	Urban	3	2.5	56	-51	-51		16	+6	+21
827	Overall accur	racy		85	+2	+4				

Table 3. Effect of location error on parameter values for user's and producer's accuracies. Column "A" is accuracy (expressed as a %) for the
spatially aligned reference and map data and column "UN" is the deviation in accuracy of the spatially unaligned reference and map data.
Negative deviations in column UN indicate that accuracy is lower in the spatially unaligned population. Classes are listed in order of decreasing
percent area based on the map classification.

836 a) Florid
---------------

837	User's Accuracy					Producer's Accuracy								
838		Pixe	<u>1</u>	Bloc	<u>k</u>	Poly	<u>gon</u>	Pixel		Bloc	<u>k</u>	Polyg	<u>gon</u>	
839	Class	А	UN	А	UN	А	UN	А	UN	А	UN	А	UN	
840	Forest	60	-22	55	-15	57	-10	86	-31	84	-24	85	-16	
841	Agriculture	99	-31	99	-25	99	-9	98	-30	98	-24	99	-14	
842	Urban	100	-35	100	-25	100	-14	95	-34	93	-24	100	-18	
843	Shrub	85	-38	82	-28	85	-27	92	-41	91	-31	97	-20	
844	Wetland	100	-20	100	-15	100	-8	55	-12	49	-7	49	-4	
845	Water	92	-10	92	-9	98	-3	100	-11	100	-12	100	-2	
846	Barren	100	-47	100	0	100	-58	100	-41	100	-50	100	-11	
847	Overall	86	-28	83	-20	85	-12							

850	b) Fayetteville
-----	-----------------

851		User	's Accurac	<u>ccuracy</u> <u>Producer's Accuracy</u>									
852		Pixe	<u>l</u>	Bloc	<u>k</u>	Poly	<u>gon</u>	Pixel		Bloc	<u>k</u>	Polyg	<u>gon</u>
853	Class	А	UN	А	UN	А	UN	А	UN	А	UN	Α	UN
854	Agriculture	83	-6	88	-5	91	-2	86	-11	85	-5	89	-7
855	Urban	85	-23	84	-14	90	-3	72	0	87	-5	90	-6
856	Forest	79	-30	68	-11	85	-14	69	-13	77	-15	95	-12
857	Wetland	65	+3	85	-12	97	-8	85	-34	65	-9	79	-8
858	Shrub	54	-21	61	-17	74	-38	68	-42	49	-13	83	-23
859	Water	61	-3	91	-9	94	-20	89	-49	77	-13	81	-17
860	Overall	77	-15	78	-8	89	-9						
861													
862													
863	c) Dare												
864		User	's Accurac	y				Produ	cer's Ac	<u>curacy</u>			
865		Pixe	<u>l</u>	Bloc	<u>k</u>	Poly	gon	Pixel		Bloc	<u>k</u>	Polyg	<u>gon</u>
866	Class	А	UN	А	UN	А	UN	А	UN	А	UN	Α	UN
867	Forest	90	-7	91	-4	98	-1	91	-9	90	-5	90	-3
868	Wetland	87	-7	88	-5	87	-1	87	-7	88	-4	95	-3
869	Shrub	71	-13	82	-12	82	-12	77	-24	77	-11	78	-1
870	Agriculture	77	-1	83	-5	84	+1	81	-9	82	+13	77	+3
871	Urban	56	-45	5	+2	5	-4	16	+23	22	+11	37	-23
872	Overall	85	-8	87	-5	89	-2						

#### 873 **Figure Captions**

874 Figure 1. Florida populations represented as difference maps for the three spatial assessment units, pixel,

875 block, and polygon. Agreement between the map and reference classifications is shown in white,

876 disagreement in black. Overall accuracy is 86% for the pixel unit, 83% for the block unit, and 85% for

877 the polygon unit.

the polygon unit.

878

879 Figure 2. Fayetteville populations represented as difference maps for the three spatial assessment units,

880 pixel, block, and polygon. Agreement between the map and reference classifications is shown in white, 881 disagreement in black. Overall accuracy is 77% for the pixel unit, 78% for the block unit, and 89% for 882

883

884 Figure 3. Dare populations represented as difference maps for the three spatial assessment units, pixel,

885 block, and polygon. Agreement between the map and reference classifications is shown in white,

886 disagreement in black. Overall accuracy is 85% for the pixel unit, 87% for the block unit, and 89% for

- 887 the polygon unit.
- 888

889 Figure 4. Effect of location error for a 3x3 pixel block assessment unit. A) Location of a 3x3 pixel block 890 assessment unit delineated on a Landsat image; this is what a photointerpreter would first look at to locate 891 an assessment unit selected by the sampling design. B) Block assessment unit of panel A classified by the 892 map being evaluated; this information would not be provided to the interpreters. C) Sample block unit 893 overlaid on a high resolution image used for determining the reference classification; the block is nearly 894 homogeneous forest. D) Same image as in C but with the block shifted one pixel down and one pixel to 895 the left to create a spatial misalignment of the map and reference locations; the reference classification for

896 the block must now address the mixed nature of the block as forest and pasture.

897

898 Figure 5. Systematic sample of points for selecting a sample of polygon assessment units. A polygon 899 would be included in the sample if a systematic sample point falls within the polygon. The probability 900 that a polygon is selected by the point sampling protocol is proportional to the area of the polygon.

- 901
- 902





905 Figure 2



908 Figure 3



## (A)





### 918 Supplemental Online Material – Figure Captions

- 919 Figure S1. Florida region target map National Land Cover Data (NLCD) 1992.
- 920 Figure S2. Fayetteville region target map National Land Cover Data (NLCD) 2001.
- 921 Figure S3. Dare region target map National Land Cover Data (NLCD) 2001.

# Appendix: Effect of spatial assessment unit on accuracy parameter values for additional (low accuracy) populations.

925 The populations created by shifting the Florida, Fayetteville, and Dare reference maps (Section 2.3) 926 represent additional example populations for examining how the values of the accuracy parameters 927 change as a function of spatial assessment unit. Although the primary purpose for creating the location-928 error impacted populations discussed in Section 2.3 was to examine how accuracy results changed for 929 different spatial units when location error was present, these location-error impacted populations provide 930 additional examples for evaluating changes in accuracy parameters resulting from use of different spatial 931 assessment units (related to Section 2.2). Table A1 provides a comparison of accuracy parameters 932 resulting from the different spatial units for a set of location-error affected populations in which accuracy 933 is lower than the populations used in Table 2. The general trends in the results observed from Table A1 934 are similar to those seen in Table 2 but the magnitude of the differences in accuracy obtained from 935 different spatial units is greater for the lower accuracy populations of Table A1. Specifically, the general 936 trend observed in Table 2 that pixel- and block-based accuracies are more similar to each other than they 937 are to polygon-based accuracy is observed in Table A1. However, the magnitude of the increase in block-938 based accuracy relative to pixel-based accuracy is much greater in Table A1 than in Table 2. The 939 polygon-based accuracies are similarly much higher than the pixel-based and block-based accuracies for 940 the low accuracy populations of Table A1. For example, overall accuracy for the polygon assessment is 941 15%, 18%, and 10% higher than the pixel assessment for the Florida, Fayetteville, and Dare populations 942 in Table A1, whereas the corresponding changes in overall accuracy are -1%, +12%, and +4% for the 943 Table 2 examples. In Table A1, block- and polygon-based user's and producer's accuracies are almost 944 always higher than the corresponding pixel-based accuracies, with increases in accuracy of 10-25% not 945 unusual. Generally the increase in accuracy of the block- and polygon-based assessments relative to the 946 pixel-based accuracy is smaller if the land-cover class comprises a relatively large percentage of the area 947 (e.g., forest in Florida and Dare, and agriculture in Fayetteville). The results of Appendix Table A1 948 suggest that differences in accuracy resulting from different spatial units are magnified when accuracy is

949 lower.

**Table A1**. Parameter values for user's and producer's accuracies resulting from different spatial

951 assessment units for three additional populations. If a block or polygon is not homogeneous in terms of

952 the map or reference classification, the mode class is used and agreement is defined as a match between

953 the mode map class and the mode reference class. Block and polygon accuracies are reported as

954 deviations from the pixel-based accuracy parameter values. Positive deviations indicate higher accuracy

955 for the block- or polygon-based assessment, and negative deviations indicate lower accuracy for the

block- or polygon-based assessment. Land-cover classes are listed in order of decreasing percent of map

- 957 area.
- 958

#### a) Florida (location-error impacted population)

960		Map	User's	s Accura	<u>cy (%)</u>	Producer's Accuracy (%)			
961	Class	Area (%)	Pixel	Block	Polygon		Pixel	Block	Polygon
962	Ag	20	68	+6	+22		68	+6	+17
963	Wetland	20	80	+5	+12		43	-1	+2
964	Forest	19	38	+2	+9		55	+5	+14
965	Urban	18	65	+20	+21		61	+8	+21
966	Shrub	16	47	+7	+11		51	+9	+26
967	Water	7	82	+1	+13		89	-1	+9
968	Barren	<1	53	+47	-11		59	-9	+30
969	Overall		58	+5	+15				

970

b) Fayetteville (location-error impacted population)

972		Map	User's	s Accura	<u>cy (%)</u>	Producer's Accuracy (%)			
973	Class	Area (%)	Pixel	Block	Polygon		Pixel	Block	Polygon
974	Ag	32	77	+6	+12		75	+5	+7
975	Urban	31	62	+8	+25		72	+10	+12
976	Forest	17	49	+8	+22		56	+6	+27
977	Wetland	10	68	+5	+21		51	+5	+20
978	Shrub	9	33	+11	+3		26	+10	+34
979	Water	<1	58	+24	+16		40	+24	+24
980	Overall		62	+8	+18				

984	4 Map			s Accura	<u>cy (%)</u>	Producer's Accuracy (%)			
985	Class	Area (%)	Pixel	Block	Polygon	Pixel	Block	Polygon	
986	Forest	45	83	+4	+14	82	+3	+5	
987	Wetland	34	80	+3	+6	80	+4	+12	
988	Ag	9	76	+2	+9	72	+3	+8	
989	Shrub	9	58	+12	+12	53	+13	+24	
990	Urban	3	11	-4	-10	39	-6	-25	
991	Overall		77	+5	+10				

983 c) Dare (location-error impacted population)