

1                   **Using Volunteered Geographic Information (VGI) in Design-Based Statistical Inference**  
2                   **for Area Estimation and Accuracy Assessment of Land Cover**

3  
4                   Stephen V. Stehman<sup>a</sup>, Cidália C. Fonte<sup>b</sup>, Giles M. Foody<sup>c</sup>, Linda See<sup>d</sup>  
5

6  
7                   <sup>a</sup> Department of Forest and Natural Resources Management, SUNY College of Environmental Science  
8                   and Forestry, Syracuse, NY 13210, United States (svstehma@syr.edu)  
9

10                  <sup>b</sup> Departamento de Matemática, Faculdade de Ciências e Tecnologia, Universidade de Coimbra, Apartado  
11                  3008, EC Santa Cruz, 3001 – 501 Coimbra, Portugal (cfonte@mat.uc.pt)  
12

13                  <sup>c</sup> School of Geography, University of Nottingham, Sir Clive Granger Building, University Park,  
14                  Nottingham, NG7 2RD, United Kingdom (giles.foody@nottingham.ac.uk)

15                  <sup>d</sup> International Institute for Applied Systems Analysis (IIASA), Schlossplatz 1, A-2361 Laxenburg, Austria  
16                  (see@iiasa.ac.at)  
17

18                  **Corresponding Author:** Stephen V. Stehman (svstehma@syr.edu)  
19

20                  **Abstract**

21                  Volunteered Geographic Information (VGI) offers a potentially inexpensive source of reference data for  
22                  estimating area and assessing map accuracy in the context of remote-sensing based land-cover  
23                  monitoring. The quality of observations from VGI and the typical lack of an underlying probability  
24                  sampling design raise concerns regarding use of VGI in widely-applied design-based statistical inference.  
25                  This article focuses on the fundamental issue of sampling design used to acquire VGI. Design-based  
26                  inference requires the sample data to be obtained via a probability sampling design. Options for  
27                  incorporating VGI within design-based inference include: 1) directing volunteers to obtain data for  
28                  locations selected by a probability sampling design; 2) treating VGI data as a “certainty stratum” and  
29                  augmenting the VGI with data obtained from a probability sample; and 3) using VGI to create an  
30                  auxiliary variable that is then used in a model-assisted estimator to reduce the standard error of an  
31                  estimate produced from a probability sample. The latter two options can be implemented using VGI

32 data that were obtained from a non-probability sampling design, but require additional sample data to  
33 be acquired via a probability sampling design. If the only data available are VGI obtained from a non-  
34 probability sample, properties of design-based inference that are ensured by probability sampling must  
35 be replaced by assumptions that may be difficult to verify. For example, pseudo-estimation weights can  
36 be constructed that mimic weights used in stratified sampling estimators. However, accuracy and area  
37 estimates produced using these pseudo-weights still require the VGI data to be representative of the full  
38 population, a property known as “external validity”. Because design-based inference requires a  
39 probability sampling design, directing volunteers to locations specified by a probability sampling design  
40 is the most straightforward option for use of VGI in design-based inference. Combining VGI from a non-  
41 probability sample with data from a probability sample using the certainty stratum approach or the  
42 model-assisted approach are viable alternatives that meet the conditions required for design-based  
43 inference and use the VGI data to advantage to reduce standard errors.

44

45 **Key Words:** probability sampling; external validity; pseudo-weights; data quality; model-based  
46 inference; Volunteered Geographic Information (VGI); crowdsourcing

47

## 48 **1. Introduction**

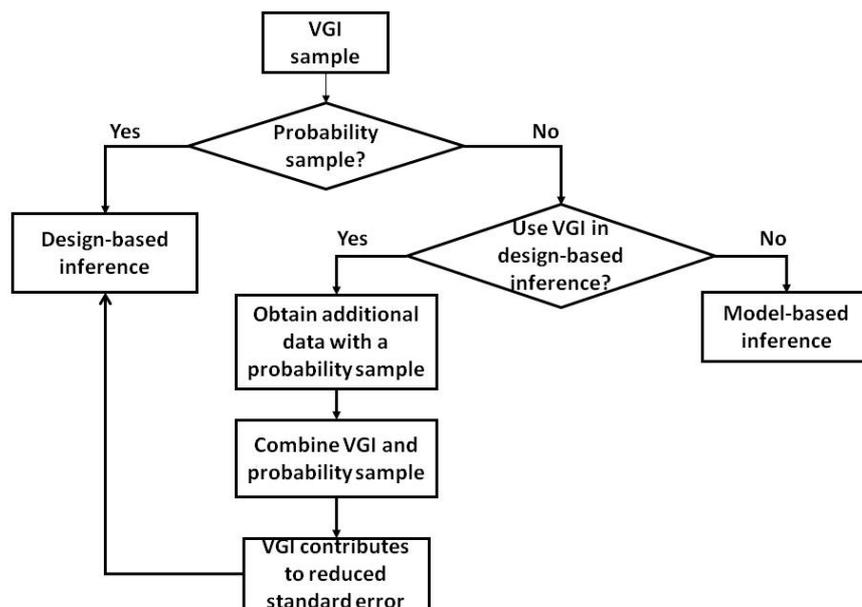
49 Volunteered Geographic Information (VGI) is defined as “tools to create, assemble, and  
50 disseminate geographic data provided voluntarily by individuals” (Goodchild 2007). For land-cover  
51 studies, VGI may provide the reference condition or the information used to determine the reference  
52 condition of a spatial unit. The reference condition, defined as the best available assessment of the  
53 ground condition, plays a critical role in accuracy assessment and area estimation (Olofsson et al. 2014).  
54 When used in map production, VGI could form all or part of the data used to train the land-cover  
55 classification algorithm. The focus of this article is the contribution of VGI to the reference data used for

56 accuracy assessment and area estimation. Accuracy assessment is an essential component of a rigorous  
57 mapping-based analysis of remotely sensed data as without it the obtained products are little more than  
58 pretty pictures and simply untested hypotheses (McRoberts 2011; Strahler et al. 2006). In addition an  
59 accuracy assessment adds value to a study, especially when estimates of class area (e.g. deforestation)  
60 are to be obtained (Olofsson et al. 2014). Fonte et al. (2015) examined the use of VGI for land cover  
61 validation, including the types of VGI that have been used, the main issues surrounding VGI quality  
62 assessment, and examples of VGI projects that have collected data for validation purposes. We build  
63 upon this past work to focus on the issue of statistical inference when incorporating VGI in applications  
64 of accuracy and area estimation, but our work is also relevant to application of citizen science data in  
65 general (Bird et al. 2014).

66 Map accuracy assessment is a spatially explicit comparison of the map class label to the  
67 reference condition on a per spatial unit basis (e.g., pixel, block, or segment). Accuracy assessment  
68 typically focuses on producing an error matrix and associated summary measures including overall,  
69 user's, and producer's accuracies (see Section 2 for details). Estimates of area of each land-cover class  
70 or type of land-cover change based on the reference condition are often produced in conjunction with  
71 the accuracy estimates (Olofsson et al. 2013, 2014). Sampling, defined as selecting a subset of the  
72 population, is almost always necessary because it is too costly to obtain a census of the reference  
73 condition. VGI represents a subset of the population and as such may be viewed as a sample. Whether  
74 the VGI data were collected via a probability sampling design is a key consideration when evaluating the  
75 utility of VGI for design-based inference. Design-based inference is a standard, widely used approach  
76 adopted in environmental science for furthering knowledge and understanding on the basis of a sample  
77 of cases rather than a study of the entire population.

78 We describe options for incorporating VGI into map accuracy assessment and area estimation  
79 within the design-based inference framework (Figure 1). We evaluate how the potential cost savings of

80 VGI can be transformed into more precise estimators (i.e., smaller standard errors, a desirable outcome  
 81 of an effective sampling strategy) within the scientifically defensible framework provided by design-  
 82 based inference. If the VGI data are obtained via a probability sampling design, application of design-  
 83 based inference is straightforward and can be informed by good practice guidelines (Olofsson et al.  
 84 2014). Alternatively, if the VGI data are not obtained via a probability sampling protocol, the VGI data  
 85 can be combined with additional data from a probability sample to produce estimates that satisfy the  
 86 conditions underlying design-based inference. In such cases the VGI data from a non-probability sample  
 87 serve as a means to reduce standard errors of estimates rather than as the sole data from which the  
 88 area and accuracy estimates are produced.



89  
 90 **Figure 1. Schema for methodologies using VGI in accuracy assessment and area estimation.**

91  
 92 This article has two major objectives. First, it illustrates how statistically rigorous and credible  
 93 inference may be drawn from studies that use VGI and thereby helps ensure that the vast potential of  
 94 VGI that has recently arisen is realized fully. This in turn will help remote sensing achieve its full

95 potential as a source of land cover information which is often constrained by lack of ground reference  
96 data. Second, the article provides methodological rigor and good practice advice for the use of data  
97 acquired via popular sample designs, ranging from judgmental to probability sampling. As such this  
98 article articulates methodology for producing credible inference from data sets that often do not  
99 conform to the requirements of widely used statistical inferential methods for two common and  
100 important application areas of remote sensing, accuracy assessment and area estimation. To do this,  
101 we, for the first time, synthesize methods developed in the general sampling literature into a  
102 comprehensive treatment of the theory and methods for using VGI in design-based inference. This  
103 includes translating methods developed for the use of non-probability samples for accuracy assessment  
104 and area estimation applications. As such we will show how VGI may be constructively used to decrease  
105 costs and reduce uncertainty (e.g., yield smaller standard errors and hence narrower confidence  
106 intervals) while following a methodology that allows for rigorous design-based inference. Throughout  
107 this article, guidance for using VGI in design-based inference is framed by examining the direct  
108 connection of the inference process to the three component protocols of accuracy assessment, the  
109 response design, sampling design, and analysis (Stehman and Czaplewski 1998).

110         The article is organized as follows. In Section 2, we define inference and describe the conditions  
111 needed to satisfy design-based inference. Considerations regarding the use of VGI in design-based  
112 inference are then explained in Section 3 in regard to the response design, sampling design and analysis  
113 protocols. Section 4 provides the details of two methods for incorporating VGI in estimation of accuracy  
114 and area that satisfy conditions of design-based inference, with both methods requiring that an  
115 additional probability sample exists or could be acquired if the VGI did not originate from a probability  
116 sampling design. Options for analysis when the only data available are VGI from a non-probability  
117 sample are discussed in Section 5. Sections 6 and 7 provide discussion and a summary of the article.

118

119

## 120 **2. Inference**

121           Following Baker et al. (2013, p.91), we define statistical inference as "... a set of procedures that  
122 produces estimates about the characteristics of a target population and provides some measure of the  
123 reliability of those estimates." Statistical inference focuses on the use of sample data to estimate  
124 parameters of a target population, where a parameter is defined as a number describing the population  
125 (e.g., the population mean and population proportion are two common parameters). Determining the  
126 numerical value of a parameter would require a census of the study region, but in practice parameters  
127 are estimated from a sample. Statistical inference also includes how bias and variance of these sample-  
128 based estimators are defined. Baker et al. (2013, p.91) further specify that "A key feature of statistical  
129 inference is that it requires some theoretical basis and explicit set of assumptions for making the  
130 estimates and for judging the accuracy of those estimates." Consequently, sampling design and analysis  
131 protocols must adhere to certain rules of implementation to ensure that the underlying mathematical  
132 basis of the inference framework is satisfied. Failure to adhere to these rules may lead to substantial  
133 bias in the estimators of parameters of interest or even nullify the ability to implement design-based  
134 inference entirely (see Section 3.3).

135           Two general types of inference are design-based inference and model-based inference (De  
136 Gruijter and Ter Braak 1990; Särndal et al. 1992; Gregoire 1998; Stehman 2000; McRoberts 2010, 2011).  
137 In design-based inference, bias and variance of an estimator are determined by the randomization  
138 distribution of the estimator which is represented by the set of all possible samples that could be  
139 selected from the population using the chosen sampling design. This randomization distribution is  
140 completely dependent on the sampling design hence the origin of the name "design-based" inference.  
141 The inclusion probabilities of the sampling design are the critical link to the randomization distribution

142 that underlies design-based inference (Särndal et al. 1992, section 2.4). The practical considerations for  
143 using VGI in design-based inference are explained in detail in Section 4.

144 A probability sampling design must satisfy two criteria related to the inclusion probabilities  
145 determined by the sample selection protocol. The inclusion probability of a particular element of the  
146 population (e.g., a pixel) is defined as the probability of that element being included in the sample. An  
147 inclusion probability is defined in the context of all possible samples that could be selected for a given  
148 sampling design. For example, if the design is simple random sampling of  $n$  elements selected from the  
149  $N$  elements of the population, the inclusion probability of each element  $u$  of the population is  $\pi_u=n/N$ .  
150 That is, in the context of all possible simple random samples of size  $n$  from this population, element  $u$   
151 has the probability of  $n/N$  of being included in the sample selected. The two requirements of a  
152 probability sampling design are that  $\pi_u$  must be known for each element of the sample and  $\pi_u>0$  for  
153 each element of the population (Särndal et al. 1992; Stehman 2000). Probability sampling requires a  
154 randomization mechanism to be present in the selection protocol. Convenience, judgment, haphazard,  
155 and purposive selection of sample elements are examples of protocols that do not satisfy the criteria  
156 defining a probability sampling design (Cochran 1977, Sec. 1.6). Use of such samples for inference  
157 carries considerable risk due to lack of representation of the population.

158 An alternative to design-based inference is model-based inference (Valliant et al. 2000). As the  
159 name implies, model-based inference requires specification of a statistical model and inference is  
160 dependent on the validity of the model. Consequently, verifying model assumptions is a critical and  
161 often challenging feature of model-based inference. Model-based inference does not require a  
162 probability sampling design, although implementation of a probability sampling design is often  
163 recommended to ensure objectivity in sample selection because of the randomization (Valliant et al.  
164 2000, p.20). Applications of model-based inference are briefly discussed in Section 5.3.

165

166

### 167 **3. Component Protocols of Accuracy Assessment and Area Estimation**

168           We describe the role of each of the three components of the methodology (response design,  
169 sampling design, and analysis) in determining how VGI can be incorporated in rigorous design-based  
170 inference. The response design is the protocol for determining the reference condition (i.e., the best  
171 available assessment of the ground condition). The response design includes all steps leading to  
172 assignment of the reference condition label of a point or spatial unit (e.g., a land-cover class or change  
173 versus no change label). The sampling design is the protocol for selecting the sample units at which the  
174 response design will be applied. Lastly, the analysis consists of defining parameters to describe  
175 properties of the population (e.g., overall accuracy, proportion of area of each class) and the formulas  
176 required to estimate these population parameters from the sample data. To justify the requirements of  
177 each step to achieve the final accuracy or area estimates, our description starts with the analysis  
178 (Section 3.1) focusing on how the VGI data would be used, followed by the steps of the response design  
179 (Section 3.2) and the sampling design (Section 3.3).

180

#### 181 **3.1 Analysis: Accuracy and Area Estimation Based on Totals**

182           The details of the analysis protocol that specify how the estimates of accuracy and area are  
183 produced yield insights into how VGI should be evaluated for use in design-based inference. The  
184 analysis focuses on summarizing information contained in an error matrix. We define the population to  
185 be a collection of  $N$  equal-area units partitioning the region of interest. The population error matrix  
186 resulting from a census can be constructed in terms of area as illustrated by the numerical example in  
187 Table 1 for a simple two-class legend, “crop” and “not crop” for a population (target region) of 1000  
188 km<sup>2</sup>. The error matrix expressed in terms of area (Table 1) could easily be converted to proportion of  
189 area by dividing each cell of the error matrix by 1000 km<sup>2</sup>. However, it is useful to focus on the error

190 matrix expressed in terms of area because we can formulate the population parameters of interest for  
 191 accuracy and area as totals or ratios of totals of areas. For example, overall accuracy is the total area of  
 192 agreement obtained from the sum of the area of the diagonal cells (930 km<sup>2</sup>) divided by the total area of  
 193 the target region (1000 km<sup>2</sup>) to yield overall accuracy of 0.93 or 93%. User's accuracy for the crop class  
 194 is the total area where both the map and reference condition are crop (840 km<sup>2</sup>) divided by the total  
 195 area mapped as crop (890 km<sup>2</sup>) to yield the parameter 0.94 or 94%. Producer's accuracy for the crop  
 196 class is the total area where both the map and reference condition are crop (840 km<sup>2</sup>) divided by the  
 197 total area of reference condition of crop (860 km<sup>2</sup>) to yield the parameter 0.98 or 98%. Lastly, the area  
 198 of reference condition of the crop class is also simply a total, in this case the sum of the two cells in the  
 199 "crop" column of reference condition (840+20 = 860 km<sup>2</sup>).

200

201 **Table 1.** Population error matrix expressed in terms of area (km<sup>2</sup>) for a hypothetical target region of  
 202 1000 km<sup>2</sup>. Overall accuracy is 93% (930/1000).

	<u>Reference Condition</u>			
<u>Map</u>	Crop	Not Crop	Total	User's
Crop	840	50	890	0.94
<u>Not Crop</u>	20	90	110	<u>0.82</u>
Total	860	140	1000	
Producer's	0.98	0.64		

209

210 Given that the parameters of interest for accuracy and area can be expressed in terms of totals,  
 211 the analysis focuses on estimating these totals. Basic sampling theory provides an unbiased estimator of  
 212 a population total in the form of the Horvitz-Thompson estimator (Horvitz and Thompson 1952). The  
 213 population total of the variable  $y_u$  is defined as

214 
$$Y = \sum_P y_u \quad [1]$$

215 where the summation is over all  $N$  elements of the population,  $P$ . For example, if  $y_u$  is the area of crop  
 216 (as determined from the reference condition) for element  $u$ , then  $Y$  is the total area of crop. The  
 217 population total  $Y$  can be estimated from a sample using the Horvitz-Thompson estimator

218 
$$\hat{Y} = \sum_s \frac{y_u}{\pi_u} \quad [2]$$

219 where the summation is over all elements of the sample  $s$ .

220 The Horvitz-Thompson estimator is an unbiased estimator of a population total for any sampling  
 221 design as long as the inclusion probabilities of the sample elements are known for that design. A useful  
 222 re-expression of the Horvitz-Thompson estimator highlighting the sample estimation weights is

223 
$$\hat{Y} = \sum_s w_u y_u \quad [3]$$

224 where  $w_u = 1/\pi_u$  is the estimation weight for element  $u$  of the sample. Because  $w_u \geq 1$ , the  $y_u$  value for  
 225 each sampled element is multiplied by an “expansion factor”  $w_u$  to estimate a total. In effect each  
 226 sample element must account for itself along with some additional elements of the population that  
 227 were not selected into the sample. For example, for simple random sampling  $w_u = N/n$  so  $y_u$  for each  
 228 sampled element is “expanded” by the multiplier  $w_u$  to account for  $N/n$  elements of the population. The  
 229 critical importance of known inclusion probabilities for rigorous design-based inference is evident via  
 230 the role of the weights  $w_u = 1/\pi_u$  in the estimator  $\hat{Y}$  (equations 2 and 3).

231 Parameters such as user’s accuracy and producer’s accuracy are ratios of totals and  
 232 consequently can be estimated by the corresponding ratio of estimated totals (Särndal et al. 1992,  
 233 section 5.3). For example, if we define  $Y$  as the total area of the population for which both the map and  
 234 reference condition are crop and  $X$  as the total area mapped as crop, the ratio of population totals  $Y/X$   
 235 would be the population parameter for user’s accuracy of crop. User’s accuracy could then be estimated  
 236 from the sample data using a ratio of Horvitz-Thompson estimators,  $\hat{Y}/\hat{X}$ , where both  $\hat{Y}$  and  $\hat{X}$  are  
 237 estimated totals based on equation (2), considering, respectively,  $y_u$ =area of pixel  $u$  with both map and

238 reference condition of crop and  $x_u$ =area of pixel  $u$  mapped as crop. In the case of a pixel-based  
239 assessment and assuming all pixels are equal area, user's accuracy of crop estimated using a ratio of  
240 Horvitz-Thompson estimators would simply require defining  $y_u=1$  if pixel  $u$  has both map and reference  
241 labels of crop ( $y_u=0$  otherwise) and defining  $x_u=1$  if pixel  $u$  has map label of crop ( $x_u=0$  otherwise). In this  
242 formulation of user's accuracy, the ratio  $Y/X$  is the proportion of pixels mapped as the target class that  
243 have the reference label of that class.

244           Formulas for the variance and estimated variance of the Horvitz-Thompson estimator are  
245 provided by Särndal et al. (1992, section 2.8). The square root of the estimated variance (standard  
246 error) would be used to construct a confidence interval for the parameter of interest so issues of  
247 inference obviously extend to variance and confidence interval estimation. Although we do not delve  
248 into the details of the formulas for variance estimators, we emphasize that known inclusion probabilities  
249 are an essential feature of variance estimation. Consequently, the requirement of implementing  
250 probability sampling to ensure known inclusion probabilities for estimating a total applies as well to  
251 estimating the variance of an accuracy or area estimator.

252           The conditions required for VGI to be used in design-based inference are apparent from the  
253 analysis protocol. The accuracy and area parameters of interest can be expressed as population totals  
254 or ratios of population totals and these totals can be estimated using the Horvitz-Thompson estimator.  
255 From the Horvitz-Thompson estimator formula (equations 2 and 3) we observe that the key features of  
256 VGI relevant to estimating a total are quality of the observation  $y_u$  and knowledge of the inclusion  
257 probability  $\pi_u$ . In other words, the questions pertinent to evaluating the utility of VGI for design-based  
258 inference are: 1) What is the quality of  $y_u$  (an issue to address in the response design) and 2) Is  $\pi_u$  known  
259 (an issue to address in the sampling design)? The following two subsections address issues of VGI  
260 related to the response and sampling designs.

261

## 262 3.2 Response Design

263 The response design is the protocol for determining the reference condition of an element of  
264 the population. In the case of a land-cover legend based on a conventional hard classification, the  
265 response design results in a reference land-cover label assigned to each pixel (i.e., if the legend consists  
266 of  $C$  classes, one and only one of these class labels is assigned to the pixel). The reference class labels  
267 can be translated to a quantity by the simple process of defining  $y_u = 1$  if pixel  $u$  has reference class  $c$  and  
268  $y_u = 0$  otherwise. Thus for example if class  $c$  is forest, all pixels with reference class forest would be  
269 assigned  $y_u = 1$  and all non-forest pixels would have  $y_u = 0$ . Evaluating and assuring the quality of VGI is  
270 critical because high quality reference data are absolutely essential to accuracy and area estimation. If  
271 the reference labels are not accurate, these errors can have a substantial impact on accuracy and area  
272 estimates (Foody 2009, 2010). Very accurate reference data obtained within a timeframe corresponding  
273 to the date of remote sensing image acquisition are a necessity for every application of accuracy  
274 assessment and area estimation from remote sensing. VGI has considerable potential as a source of  
275 reference data, notably in facilitating the collection of a large set of observations over broad  
276 geographical regions. However, the use of volunteers rather than experts in assigning the reference  
277 class labels may exacerbate concerns regarding label accuracy, although amateurs can sometimes be as  
278 accurate as experts in labeling (See et al. 2013). Further, VGI tends to be collected continuously rather  
279 than within a narrow time frame which can limit its value, especially for studies of land-cover change.

280 Applications in which VGI has been collected for land cover and land use studies are becoming  
281 increasingly common. Fonte et al. (2015) reviewed several applications including:

- 282 1) Geo-Wiki project, which uses the crowd for interpretation of very high resolution satellite  
283 imagery (Fritz et al. 2012);
- 284 2) VIEW-IT, which is a validation system for MODIS land cover (Clark and Aide 2011); and

285 3) geo-tagged photographs for land cover validation from different applications such as the  
286 Degree Confluence Project, Geograph, Panoramio and Flickr (Antoniou et al. 2016; Fonte et al.  
287 2015; Iwao et al. 2006).

288 Another source of VGI for land-cover studies is the LACO-Wiki system, an online land cover validation  
289 tool intended as a repository of openly available validation data crowdsourced from different users (See  
290 et al. 2017). More recently, land cover and land use have been crowdsourced in the field through the  
291 FotoQuest Austria app, which sends users to specific locations and loosely follows the LUCAS protocol  
292 for data collection (Laso Bayas et al. 2017). Hou et al. (2015) describe geo-tagged web texts as an  
293 alternative to photographs as yet another source of VGI useful for land-cover studies.

294 The quality of the VGI data collected for land cover and land use studies has received recent  
295 attention. A substantial body of literature focuses on the positional quality and completeness of  
296 OpenStreetMap (OSM), the most commonly cited VGI project (e.g., Ciepluch et al. 2010; Girres and  
297 Touya 2010; Haklay 2010). Other elements of quality include thematic accuracy (which is relevant to  
298 land cover and land use), temporal quality, logical consistency, and usability, all of which are set out in  
299 ISO 19157 (Fonte et al. 2017a). In addition, Antoniou and Skopeliti (2015) outline quality indicators that  
300 are tailored to VGI such as data indicators, demographic and other socio-economic indicators, and  
301 indicators about the volunteers. Due to the specificities of VGI when compared to traditional  
302 geographic information and the diversity of uses of these data, additional methodologies are starting to  
303 be developed that aim to integrate several quality measures and indicators into quality assessment  
304 workflows, enabling quality data to be combined to produce more reliable quality information (e.g.,  
305 Bishr and Mantelas 2008; Jokar Arsanjani and Bakillah 2015; Meek et al. 2016).

306 Although concern with reference data error may be heightened when VGI is used, there are  
307 methods such as latent class analysis, which can be used to characterize volunteers in terms of their  
308 quality in labeling classes and could therefore be used to filter or weight the data when used

309 subsequently in applications (Foody et al. 2013, 2015). These issues of data quality associated with the  
310 response design are critical to the overall process of accuracy and area estimation. In reality, reference  
311 data quality issues are equally impactful whether the source of the reference classification is VGI or  
312 expert interpretation (See et al. 2013).

313

### 314 **3.3 Sampling Design**

315 The sampling design is the protocol used to select the subset of locations (e.g., pixels) at which  
316 the reference condition is determined. As noted earlier, the inclusion probability of pixel  $u$  is denoted as  
317  $\pi_u$ , and the two criteria defining a probability sampling design are: 1)  $\pi_u$  is known for all pixels in the  
318 sample and 2)  $\pi_u > 0$  for all pixels in the population. Because probability sampling is a requirement of  
319 rigorous design-based inference, the sample selection protocol must ensure that these two conditions  
320 of  $\pi_u$  are satisfied. Moreover, randomization of the sample selection is required of all probability  
321 sampling designs as it is this randomization that creates the probabilistic foundation for design-based  
322 inference. The sampling design is linked to the analysis via the inclusion probabilities that are  
323 incorporated in the Horvitz-Thompson estimator (equations 2 and 3).

324 Because design-based inference requires known inclusion probabilities, it is critical to establish  
325 whether a probability sampling design was the basis for collecting VGI data. The distinction between  
326 active and passive VGI is relevant in this regard. Active VGI refers to directing volunteers to specific  
327 sample locations (e.g., See et al. 2016) and therefore allows for implementing a probability sampling  
328 design for collecting VGI. Conversely, passive VGI refers to allowing volunteers to choose where they  
329 will collect data and typically leads to purposive or convenience sampling with attendant concern  
330 regarding lack of representation of the full population. The protocols that determine where VGI data  
331 are collected span a continuum ranging from rigorous probability sampling to selection by judgment or  
332 convenience without an underlying random mechanism.

333           The Degree Confluence Project (Iwao et al. 2006) is an example in which VGI data are collected  
334 via a probability sampling protocol. These data are obtained at locations defined by the intersection of  
335 lines of latitude and longitude and therefore originate from a design akin to systematic sampling (due to  
336 the Earth's shape the distances between sample points vary with latitude so the inclusion probabilities  
337 would not all be equal but would still be known). A second example of VGI based on a probability  
338 sampling design is the FotoQuest Austria app which uses the Land Use/Cover Area frame Survey (LUCAS)  
339 sample (which is based on a systematic sample of points spaced 2 km apart in the four cardinal  
340 directions across the European Union) followed by a stratified sample (Martino et al. 2009). That is, land  
341 cover and land use were crowdsourced via the FotoQuest Go mobile app in which volunteers were sent  
342 to specific locations that formed part of the LUCAS systematic sample for Austria, and the LUCAS sample  
343 was then augmented with additional sample units (Laso Bayas et al. 2016).

344           Several VGI applications include sample data originating from both probability sampling designs  
345 and volunteer chosen locations. The Geo-Wiki project is used to collect land cover and land use data via  
346 different campaigns (See et al. 2015). These campaigns have all had different purposes and hence were  
347 driven by different sampling designs. For example, the first campaign to validate a map of land  
348 availability for biofuels was driven by a stratified random sample with equal sample size in both the land  
349 available stratum and the land unavailable stratum. To this an additional sample from cropland areas  
350 was added although the data were not used to undertake an accuracy assessment as such but to modify  
351 the statistics on how much land is available (Fritz et al. 2013). Other studies have made use of Geo-Wiki  
352 data from previous campaigns for validation that were not obtained using a probability sampling  
353 approach for the specific product to be validated (see, for example, Schepaschenko et al. (2015) and  
354 Tsendbazar et al. (2015) for review of reference datasets including those from Geo-Wiki). The VIEW-IT  
355 application (Clarke and Aide 2011) either directs users to specific locations selected based on a  
356 probability sampling design or users can provide information about the land cover at any location, which

357 means these latter sample locations would not be part of a probability sampling design. The LACO-Wiki  
358 system (See et al. 2017) has built-in probability sampling schemes although users can upload their own  
359 sample locations that do not necessarily conform to a probability sampling design.

360           Photograph repositories such as Panoramio, Flickr, and Instagram are examples of passive VGI  
361 and therefore do not conform to any probability sampling design. For example, photographs made  
362 available by citizens may be positioned at any location chosen by the volunteer (such as the  
363 photographs available in Flickr or Instagram), or collected at predefined locations. Similarly, the data  
364 available in collaborative projects such as OSM are created at locations of interest to the citizen  
365 volunteers, and consequently these data have no underlying probability sampling design. The amount  
366 and quality of the OSM data are known to be correlated with demographic or socio-economic factors  
367 (e.g., Mullen et al. 2014; Elwood et al. 2013) and this offers some possibility for adjusting estimates to  
368 account for misrepresentation of the population (see Section 5.1).

369           The Geograph project asks users to take photographs in every square kilometer of the United  
370 Kingdom and classify them (now also extended to other locations in the world). Since 2005, 83.4% of  
371 the 1 km<sup>2</sup> squares in Great Britain and Ireland have photographs (<http://www.geograph.org.uk/>,  
372 accessed 29 October 2017) and nearly 5.5 million images are available within this time period.  
373 Volunteers may choose locations within each square kilometer at which photographs are taken.  
374 Therefore, if each photograph is viewed as representing a point location or, for example, the 30 m x 30  
375 m pixel surrounding the photograph's location, the data would not meet the criteria defining a  
376 probability sampling design due to the lack of randomization in the selection protocol. Directing the  
377 volunteers to cover the 1 km<sup>2</sup> squares provides a better degree of spatial representation of the VGI than  
378 might otherwise occur if volunteers are allowed to choose locations completely on their own.  
379 Specifically, the 1 km<sup>2</sup> squares effectively serve as spatial (geographic) strata, and with over 83% of  
380 these strata visited, the Geograph project data achieve the desirable design criterion of being spatially

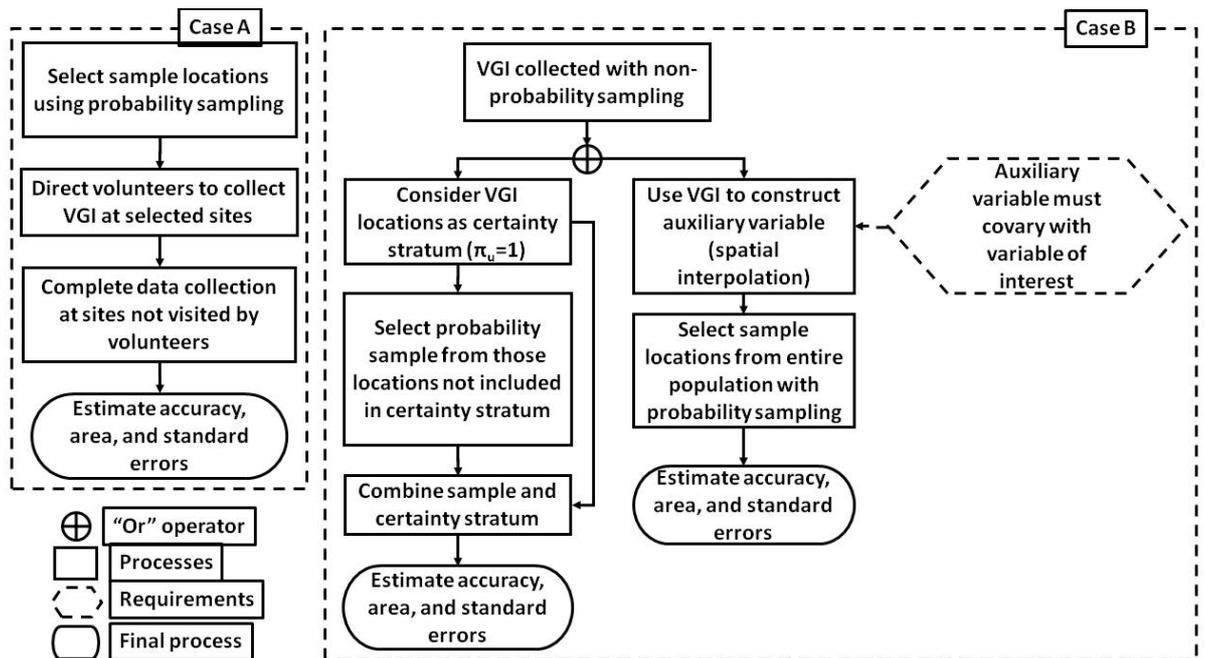
381 well distributed (Stehman 1999, Figure 3). The Geograph project data collection protocol illustrates the  
382 fact that within the class of non-probability sample designs, features can be built into the protocol to  
383 enhance representation of the VGI data.

384

#### 385 **4. Methods to Use VGI in Design-based Inference**

386 In this section, we address how to incorporate VGI into design-based inference focusing on  
387 sampling design and estimation considerations (Figure 2). The label quality issues of VGI remain a  
388 concern but are not addressed in this section. The most straightforward approach to ensure the utility  
389 of VGI for design-based inference is to direct volunteers to collect data at locations specified by a  
390 probability sampling design (which is possible with “active VGI”). Several examples of VGI collections  
391 based on a probability sampling design were documented in Section 3.3. Specifying sample locations  
392 selected via probability sampling has the potential drawback that volunteer participation may be  
393 reduced if volunteers are unable to choose locations of personal interest. Consequently, additional  
394 effort may be necessary to obtain  $y_u$  at those locations neglected by volunteers.

395



396

397

**Figure 2. Schema for using VGI in design-based inference.**

398

If a large quantity of VGI obtained from a non-probability sampling design exists, the VGI data

399

may be augmented with data from a probability sampling design (Figure 2). Two options are described

400

in the following subsections. In the first option, the VGI data are treated as a “certainty stratum” and

401

combined with data from a probability sample selected from the locations not already included in the

402

VGI data. In the second option, the probability sample is selected from the full population and the VGI

403

data are used to construct an auxiliary variable that is then incorporated in a model-assisted estimator

404

to reduce the standard errors of the estimates based on the data from the probability sample.

405

#### 406 **4.1 VGI Incorporated as a Certainty Stratum**

407

VGI data can be combined with data obtained from a probability sample by treating each VGI

408

sample unit (e.g., a pixel) as belonging to a “certainty stratum” in which the inclusion probability is  $\pi_u=1$

409

(Overton et al. 1993). By assigning  $\pi_u=1$  to each VGI sample unit, we acknowledge that these sample

410

units were not selected via a randomized selection protocol, and instead we view these units as having

411 been purposely selected to be included with certainty in the sample. From the remaining units of the  
 412 population not included in the VGI certainty stratum, a probability sampling design is implemented and  
 413 these newly selected sample units are combined with the VGI data to produce the accuracy and area  
 414 estimates. In this approach the VGI data are used directly in the estimation of accuracy and area, so the  
 415 quality of the VGI data is a critical concern.

416 All sample units selected via the probability sampling design will have a known inclusion  
 417 probability and the data from these sample units can be combined with the VGI data using the Horvitz-  
 418 Thompson estimator. Specifically, suppose there are  $N_1$  elements for which we have no VGI and  $N_2$   
 419 elements for which VGI provides  $y_u$  ( $N=N_1+N_2$ ). Further, let  $G$  denote the subset for which VGI is  
 420 available (the “G” is from the middle letter of VGI) and  $\tilde{G}$  denote the subset of the population for which  
 421 VGI is not available. The population total  $Y$  can then be partitioned into summations over the two  
 422 subpopulations  $\tilde{G}$  and  $G$ ,

$$423 \quad Y = \sum_{\tilde{G}} y_u + \sum_G y_u = Y_{\tilde{G}} + Y_G \quad [4]$$

424 Because  $Y_G$  (total of  $y_u$  for the VGI data) is known, it is only necessary to estimate  $Y_{\tilde{G}}$  from the sample.  
 425 Therefore, an estimator of  $Y$  can be expressed as

$$426 \quad \hat{Y} = \sum_s y_u / \pi_u + \sum_G y_u = \hat{Y}_{\tilde{G}} + Y_G \quad [5]$$

427 where the first summation is over the elements selected in the sample from the  $N_1$  elements of the  
 428 population  $\tilde{G}$  for which VGI is not available. The variance of  $\hat{Y}$  is  $V(\hat{Y}) = V(\hat{Y}_{\tilde{G}})$  because the total of the  
 429 VGI data is a known quantity with no uncertainty attributable to sampling. That is, the only uncertainty  
 430 attributable to sampling arises from estimating the total  $Y_{\tilde{G}}$  for the non-VGI portion of the population,  
 431  $\tilde{G}$ .

432 The benefit of the VGI data when incorporated as a certainty stratum is to reduce the standard  
 433 errors of the accuracy and area estimators and accordingly to decrease the width of confidence intervals  
 434 for the parameters of interest. To illustrate the potential reduction in standard error, we focus on the

435 objective of estimating area based on the reference condition obtained for each sample unit. The  
 436 benefit of the VGI data can then be quantified by comparing the variance of the estimator of total area  
 437 without using VGI data to the variance of the estimator using the certainty stratum approach (equation  
 438 5). Several conditions are imposed to simplify the variance comparison: 1) the sample of non-VGI units  
 439 is selected by simple random sampling; 2) the VGI data have the same variability as the non-VGI data  
 440 (i.e., the variance of  $y_u$  for the VGI subpopulation  $G$  is the same as the variance of  $y_u$  for the non-VGI  
 441 subpopulation  $\tilde{G}$ ); and 3) the sample size  $n$  is the same regardless of whether VGI is present (i.e., the VGI  
 442 data are viewed as obtained at no cost so  $n$  is the same with or without VGI). If no VGI data are  
 443 available and a simple random sample is selected from the full population of  $N$  elements (i.e.,  $N_2=0$   
 444 because no VGI data exist), the variance of the estimated total is

$$445 \quad V(\hat{Y}) = N^2 \left(1 - \frac{n}{N}\right) V_y/n \quad [6]$$

446 The variance of  $\hat{Y}$  when VGI is available for  $N_2$  elements of the subpopulation  $G$  is derived as follows. A  
 447 simple random sample of  $n$  elements is selected from the  $N_1$  non-VGI units. The variance of the  
 448 estimated total combining the VGI data with the non-VGI sample (equation 5) depends only on the  
 449 variance of the total estimated from the non-VGI sample units,

$$450 \quad V(\hat{Y}_{\tilde{G}}) = N_1^2 \left(1 - \frac{n}{N_1}\right) V_y/n \quad [7]$$

451 To quantify the reduction in variance achieved by the VGI data, we examine the ratio of the two  
 452 variances,

$$453 \quad R = \frac{V(\hat{Y}_{\tilde{G}})}{V(\hat{Y})} = \frac{N_1^2 \left(1 - \frac{n}{N_1}\right)}{N^2 \left(1 - \frac{n}{N}\right)} \quad [8]$$

454 The  $V_y/n$  term common to both equations (6) and (7) cancels in the ratio  $R$  by virtue of the assumption  
 455 that the variability of  $y_u$  is the same in the VGI and non-VGI subpopulations (if  $V_y$  is different in the two  
 456 subpopulations,  $R$  will be impacted by the ratio of the variances of the two subpopulations,  $G$  and  $\tilde{G}$ ).

457 Under the assumption of equal variance for the two subpopulations, the benefit of VGI to  
 458 reduce variance depends on the proportion of the population that is covered by the VGI data, which is  
 459 defined as  $k=N_2/N$ . If we define  $f=n/N$  to be the proportion of the total population selected for the  
 460 probability sample, then  $R$  can be re-written as

$$461 \quad R = (1 - k)(1 - f - k)/(1 - f). \quad [9]$$

462 If no VGI data exist, then  $k=0$  and  $R=1$  as expected because there would be no reduction in variance  
 463 from VGI. Conversely, if  $k=1$ , then  $R=0$  as expected because the VGI would constitute a census and the  
 464 population total  $Y$  would be known yielding a variance of 0. As the quantity of VGI gets larger (i.e.,  
 465  $k=N_2/N$  increases),  $R$  decreases indicating a greater benefit accruing to the availability of the VGI data.  
 466 Numerical values of  $\sqrt{R}$  (ratio of standard errors) for several combinations of  $k$  and  $f$  are presented in  
 467 Table 2. For a fixed value of  $f=n/N$ ,  $\sqrt{R}$  decreases approximately linearly with increasing  $k$ . For a fixed  
 468 value of  $k$ , the decrease in  $\sqrt{R}$  is much less prominent as  $f$  increases except for the case with  $f=0.25$  and  
 469  $k=0.75$  which represents a census so  $V(\hat{Y}_{\bar{c}}) = 0$ . To simplify the problem still further, assume that the  
 470 spatial unit of the assessment is a pixel and that  $N$  is so large that  $f = n/N = 0$ . Then setting  $f = 0$  in  
 471 equation (9), we obtain  $R = (1 - k)^2$  which leads directly to

$$472 \quad \sqrt{R} = 1 - k \quad [10]$$

473 Thus for very large populations the reduction in standard error achieved by VGI will be directly related  
 474 to  $k$ , the proportion of the population for which VGI is available – the greater the quantity of VGI  
 475 available (i.e., larger  $k$ ) the greater the reduction in standard error.

476

477

478

479

480 **Table 2.** Reduction in standard error achieved by using VGI in the certainty stratum approach. Values  
 481 shown in the table are  $\sqrt{R}$  where R is the ratio of the variance of the estimated total with VGI data  
 482 incorporated in a certainty stratum divided by the variance of the estimated total in the absence of VGI  
 483 (see equations 8 and 9). Ratios are provided for different combinations of  $k=N_2/N$  (the proportion of the  
 484 region of interest covered by VGI) and  $f=n/N$  (proportion of the study region covered by the simple  
 485 random sample).

	$f = n/N$					
$k$	0.00	0.01	0.05	0.10	0.25	
0.01	0.99	0.99	0.99	0.99	0.99	0.99
0.05	0.95	0.95	0.95	0.95	0.95	0.94
0.10	0.90	0.90	0.90	0.89	0.88	
0.25	0.75	0.75	0.74	0.74	0.71	
0.50	0.50	0.50	0.49	0.47	0.41	
0.75	0.25	0.25	0.23	0.20	0.00	
0.90	0.10	0.10	0.07	0.00	0.00	

495  
 496 Equation (9) and the results of Table 2 can be used to examine the benefit of VGI arising from  
 497 photographs contributed by volunteers (Antoniou et al. 2016), a common source of VGI for land-cover  
 498 studies. Suppose we assume a photograph to be representative of a 30 m x 30 m pixel and consider a  
 499 region of interest that covers 8 million km<sup>2</sup> (roughly the size of the conterminous United States,  
 500 excluding Alaska and Hawaii). This region would have approximately  $N = 9$  billion pixels. To achieve a  
 501 5% reduction in the standard error of the estimated area of a targeted class (i.e.,  $\sqrt{R}$  changes from 1 to  
 502 0.95) the certainty stratum approach would require  $k=N_2/N=0.05$  which translates to needing  $N_2 = 450$   
 503 million photographs. As a second example, suppose the target region of interest covers 100,000 km<sup>2</sup>  
 504 (area slightly larger than Portugal). This population would have  $N = 100$  million pixels (30 m x 30 m) so

505 for VGI data to contribute a 5% reduction in standard error we would need  $N_2 = 5$  million photographs.  
506 Typically the VGI photographs will have to be processed to obtain the land-cover information of interest  
507 (e.g., a land-cover class). Consequently, the large number of photographs needed in these examples to  
508 achieve only a 5% reduction in standard error would require substantial computer processing capability  
509 and possibly automated methods to identify the land-cover class from the photographs. Accordingly,  
510 the response design effort to process such large numbers of photographs may make this use of VGI cost  
511 prohibitive in some applications.

512 The certainty stratum approach may have greater utility when the VGI data are in the form of  
513 fully mapped areas classified to a land-cover or change type (i.e., in contrast to individual, unlabeled  
514 photographs as in the previous paragraph). For example, Fonte et al. (2017b) described an application  
515 in which OSM provided land-cover information for two study areas of 100 km<sup>2</sup> in London and Paris.  
516 OSM coverage was 88% for the London region and 97% for the Paris region. Because of the substantial  
517 portion of area covered by OSM ( $k=0.88$  for London and  $k=0.97$  for Paris) a large reduction in standard  
518 error of accuracy and area estimates would be expected by using these OSM data in the certainty  
519 stratum approach. For example, if  $k=0.88$  and  $f=0.1$  (the London example), we obtain  $R=0.00266$   
520 ( $\sqrt{R}=0.05$ ) indicating that the standard error of the certainty stratum estimator would be 5% of the  
521 standard error of the estimated area when not using the VGI from OSM. Obviously the areas of the  
522 regions of interest for the OSM examples in this paragraph are much smaller than for the examples in  
523 the previous paragraph and  $k$  would surely be smaller if OSM were to be used for national estimates.

524

#### 525 **4.2 Use of VGI in a Model-Assisted Estimator**

526 Brus and de Gruijter (2003) developed an approach to use data from a non-probability sampling  
527 design to produce estimates within the design-based inference framework. In this approach, a spatial  
528 interpolation method is applied to the non-probability sample of VGI data to construct an auxiliary

529 variable for all  $N$  elements of the population. The auxiliary variable is then used in a model-assisted  
530 estimator to achieve a reduction in standard error. Model-assisted estimators represent a broad class of  
531 estimators in which one or more auxiliary variables are incorporated in the estimator. Common  
532 examples of model-assisted estimators include difference, ratio, and regression estimators as well as  
533 post-stratified estimators (Särndal et al. 1992; Gallego 2004; Stehman 2009; McRoberts 2011; Sannier et  
534 al. 2014). The auxiliary variables are expected to covary with the target variable of interest and the  
535 information in the auxiliary variables, when incorporated in the model-assisted estimator, thus serves to  
536 reduce standard errors (Särndal et al. 1992, Chapter 6).

537         The Brus and de Gruijter (2003) approach could be applied to VGI as follows. Consider the  
538 objective of estimating the proportion of area of a class (e.g., area of forest) based on the reference  
539 condition. Suppose the spatial unit of the analysis is a pixel and the VGI data consist of  $N_2$  pixels labeled  
540 as forest or non-forest. The Brus and de Gruijter (2003) approach uses these VGI data to construct an  
541 auxiliary variable  $x_u$  for all  $N$  pixels in the population. For example, for a binary classification of forest /  
542 non-forest, the auxiliary variable would be defined as  $x_u=1$  if the class is forest and  $x_u=0$  if the class is  
543 non-forest. The auxiliary variable  $x_u$  is known for the  $N_2$  pixels comprising the VGI, and the Brus and de  
544 Gruijter (2003) approach would then implement a spatial interpolation method such as indicator kriging  
545 (e.g., Isaaks and Srivastava 1989) to predict values of  $x_u$  for the  $N-N_2$  pixels not included in the VGI subset  
546 of the population. The binary forest / non-forest classification of the region predicted from the VGI data  
547 could be used in the same manner as auxiliary data from any forest / non-forest map. For example, to  
548 estimate the proportion of area of forest based on the reference condition ( $y_u$ ), a probability sample  
549 from all  $N$  pixels would be selected for which the reference class of each sampled pixel would be  
550 obtained. If the reference observation is also a binary forest / non-forest classification (i.e.,  $y_u=1$  if the  
551 reference condition is forest,  $y_u=0$  otherwise), an error matrix could be estimated from the sample  
552 based on the reference class data and the map classification of forest or non-forest created from the VGI

553 data. The error matrix information could then be combined with the VGI generated forest / non-forest  
554 map information to produce a post-stratified estimator of the proportion of area (Card 1982; Stehman  
555 2013). The expectation is that the auxiliary variable created from the VGI would yield a reduction in  
556 standard error of the post-stratified estimator relative to an estimator that did not incorporate the VGI.  
557 That is, the map generated via spatial interpolation of the VGI data would be used in the same way that  
558 a forest / non-forest map derived from remotely sensed data would be used in a post-stratified  
559 estimator.

560         The Brus and de Gruijter (2003) method requires a probability sample to provide the reference  
561 data ( $y_u$ ) for the accuracy and area estimates. This probability sample must be selected from the full  
562 population of  $N$  units, including those units for which VGI is available. In contrast, the certainty stratum  
563 use of VGI (section 4.1) does not require a sample from the subpopulation  $G$  that has VGI. The Brus and  
564 de Gruijter (2003) approach does not use the VGI data as the observed response (i.e., the reference data  
565 value,  $y_u$ ) so the quality of the class labels associated with the VGI data will not impact the estimates in  
566 terms of potential bias attributable to labeling error of the VGI. However, better quality (i.e., more  
567 accurate) VGI data would likely yield a greater reduction in standard error in the same manner that a  
568 more accurate map yields a greater reduction in standard error when the map data are used in a post-  
569 stratified estimator (Stehman 2013). In the context of land-cover accuracy and area estimation  
570 applications, remote sensing information is almost always available to produce a map that would  
571 provide auxiliary information that could be used in a model-assisted estimator. Spatial interpolation of  
572 VGI using the methods described by Brus and de Gruijter (2003) provides another option for producing a  
573 map of auxiliary information, and incorporating remote sensing imagery in linear spatial models (Diggle  
574 et al. 1998) might further enhance the precision benefit of the Brus and de Gruijter (2003) approach.

575         To summarize, the model-assisted estimator based on spatially interpolated data does not rely  
576 on the VGI data to provide the  $y_u$  values that are the basis of the parameter estimates thus decreasing

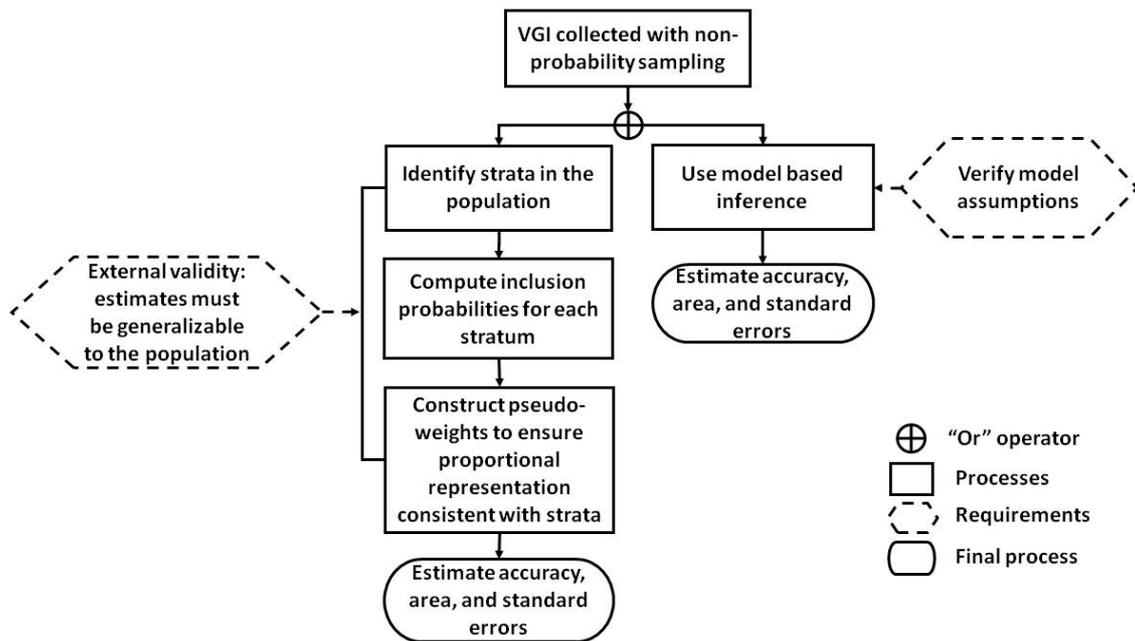
577 the concern with bias attributable to inaccurately labeled VGI data. Instead, the approach employs the  
578 VGI to create an auxiliary variable  $x_u$  that is then used in a model-assisted estimator to reduce the  
579 standard errors of the accuracy and area estimates. The magnitude of the reduction in standard error  
580 would depend on the quality of the VGI. While this approach would have great utility if no other  
581 auxiliary information were available, we typically have access to remotely sensed data that could be  
582 used to produce a classification that would serve the same purpose as a map derived from spatially  
583 interpolating VGI data. Consequently, for land-cover studies the primary benefit obtained by spatial  
584 interpolation of VGI may occur in circumstances where a map produced from remotely sensed data is  
585 not available.

586

## 587 **5. Use of VGI from Non-Probability Samples**

588 If the VGI data are the only source of reference data (i.e., there is no probability sample and  
589 unable to acquire one), it will be challenging to use these VGI data in the manner of design-based  
590 inference (Figure 3). One option for using VGI in this context is to replace the estimation weights  
591  $w_u=1/\pi_u$  (equation 3) by pseudo weights that depend on assuming the sample can be treated as though  
592 it had been obtained via a probability sampling design. For example, suppose the reference data for  
593 accuracy assessment and area estimation are land-cover interpretations extracted from a non-  
594 probability sample of photographs. If the inclusion probabilities ( $\pi_u$ ) of the spatial units represented by  
595 these photographs are unknown, one approach to estimate totals is to assume that the VGI locations  
596 represent a stratified random sample (see Section 5.1 for details). Using this approach it is possible to  
597 construct pseudo-weights such that estimated totals will match known parameters of the population.  
598 Although this weighted estimation approach can adjust a VGI sample to achieve estimates that  
599 correspond to the correct proportional representation of the population, the question of “external  
600 validity” of the VGI data must be addressed. External validity is defined and applied in Section 5.2.

601 Model-based inference is a second option for using VGI data that were not obtained from a probability  
 602 sampling design. The application of model-based inference to accuracy and area estimation is discussed  
 603 in Section 5.3.  
 604



605  
 606 **Figure 3. Schema for using VGI collected via a non-probability sampling design.**

607 **5.1 Estimation Based on Pseudo-Weights**

608 If the only reference data available for accuracy and area estimation are VGI that did not originate  
 609 from a probability sampling design, an obvious initial step in the analysis is to examine the proportional  
 610 distribution of the VGI sample relative to known characteristics of the population. For example, using a  
 611 land-cover map of the study region, we could compare the proportion of the VGI data found within each  
 612 land-cover class to the proportion of each class in the entire population. For the hypothetical numerical  
 613 example of Table 3, the VGI sample shows preferential selection from the developed and crop classes at  
 614 the expense of representation of the “other” and natural vegetation classes reflecting the relative ease  
 615 of access to the classes associated with the transport network. Representativeness of the VGI data

616 could also be assessed by examining the distribution of distances to the nearest road or distances to the  
 617 nearest population center. For example, we could compare the mean distance to the nearest road for  
 618 the VGI locations to the mean distance for all  $N$  pixels in the population. If the mean for the VGI  
 619 locations was less than the mean for the population, this discrepancy would indicate preferential  
 620 selection of VGI closer to a road. A relevant question is then whether this preferential selection could  
 621 introduce bias because map accuracy may differ depending on proximity to a road.

622

623 **Table 3.** Hypothetical data illustrating evaluation of the proportional representation of VGI. The  
 624 distribution of the percent area of the map classes is compared between the VGI sample ( $n=100$ ) and  
 625 the population (i.e., entire region) known from a land-cover map of the study region.

626

627	628 <u>Map Class</u>	627 <u>Area (%)</u>	
		VGI	Population
629	Developed	25	10
630	Crop	35	20
631	Natural vegetation	30	50
632	Other	10	20

633

634 In general, we could attempt to adjust estimates to account for recognized non-proportionality of  
 635 the VGI data relative to known population characteristics (Dever et al. 2008). For the example data of  
 636 Table 3, the difference between the distribution of the VGI and population data suggests that weighting  
 637 the data to adjust for this discrepancy would be a good idea when producing estimates. One approach  
 638 would be to construct weights such that the estimates based on the weighted analysis of the VGI data  
 639 correspond to known population quantities. A simple way to achieve this is to treat the non-probability

640 sample as having arisen from a stratified design (e.g., Loosveldt and Sonck 2008). Inclusion probabilities  
 641 for each stratum are then defined as  $\pi_u = n_h/N_h$  where  $n_h$  is the observed sample size (from the VGI  
 642 sample) in stratum  $h$  and  $N_h$  is the population size in stratum  $h$ . The estimation weight for pixel  $u$  is then  
 643  $w_u = 1/\pi_u$ , and these weights could be used in the Horvitz-Thompson estimator. These stratified  
 644 estimation pseudo-weights for the hypothetical data of Table 3 are presented in Table 4. Referring to  
 645 weights constructed in this manner as “pseudo-weights” highlights the fact that they are not derived  
 646 from inclusion probabilities generated by a probability sampling protocol.

647

648 **Table 4.** Pseudo-weights for VGI sample units based on distributions by class shown in Table 3 ( $n_h$  and  
 649  $N_h$  represent the number of pixels for each class in the VGI sample and in the population).

650

	$n_h$	$N_h$	
<u>Class</u>	<u>VGI</u>	<u>Map</u>	<u><math>w_u = N_h/n_h</math></u>
Developed	25	1000	40
Cultivated	35	2000	57
Natural veg	30	5000	167
<u>Other</u>	<u>10</u>	<u>2000</u>	<u>200</u>
Total	100	10000	

658

659 To illustrate how the stratified estimation approach using pseudo-weights is implemented, consider  
 660 estimating the proportion of area mapped as the developed class. From Table 3, we know this  
 661 proportion is 0.10 because we have the map for the entire population. How well does the VGI sample  
 662 estimate this parameter? We observe that 25 out of 100 VGI pixels are mapped as developed so the  
 663 estimated proportion of mapped developed is then 0.25 from the VGI data, greater than the known

664 parameter of 0.10 for the population. To produce the estimator using the stratified pseudo-weights of  
665 Table 4 we define  $y_u=1$  if the sample pixel has the map label of developed and  $y_u=0$  otherwise. Then for  
666 the developed class stratum,  $y_u=1$  for all 25 sample pixels and each of these pixels has a weight of  
667  $w_u=40$ , so the estimated total contributed from this stratum is  $40 \times 25 = 1,000$  pixels (using equation 3).  
668 For the other three strata,  $y_u=0$  for all sample pixels so these strata contribute no additional pixels to the  
669 estimated number of mapped developed pixels. Dividing the estimated total number of map pixels  
670 labeled as developed (1,000) by the number of pixels in the population ( $N=10,000$ ) yields an estimated  
671 proportion of 0.10 which matches the population proportion of mapped developed area from Table 3.  
672 Thus the sample estimate using the pseudo-weights matches this known population proportion.

673 In general, the pseudo-weights can be constructed so that the sample estimates will equal known  
674 population values. In the example of Table 4, the pseudo-weights reproduce the known values  
675  $N_h$ =population size of each stratum, a property known as “proportional representation.” These same  
676 estimation pseudo-weights are then applied to estimate the target population parameters and the  
677 assumption is that estimation weights that effectively adjust the VGI sample data to match known  
678 population parameters will also work well when estimating the target parameters for which we do not  
679 have full population information. Other more complex methods for creating estimation weights include  
680 raking, general calibration estimators (Deville and Särndal 1992), and propensity scores (Valliant and  
681 Dever 2011). Models can be used to produce the pseudo-weights used in lieu of weights that are the  
682 inverse of the inclusion probabilities of a probability sampling design, but Valliant (2013, p.108) points  
683 out that this approach has not yielded promising results because the models are weak and the  
684 requirements excessive for covariates to be used in the models.

685

## 686 **5.2 External validity**

687 Pseudo-estimation weights can be used to produce estimates that capture the proportional  
688 distribution of known population characteristics (i.e., covariates). However, another important aspect of  
689 representativeness of non-probability sample data is external validity, defined as the parameter estimates  
690 being “generalizable outside the sample, say to the population of interest” (Dever and Valliant 2014). For  
691 the pseudo-weight estimation approach described in the previous section, establishing external validity  
692 would require that accuracy for the subset of the population represented by the VGI locations be  
693 equivalent to accuracy of the full region. Proportional representation of the estimates (Table 4) produced  
694 from non-probability sample data is one aspect of external validity, but proportional representation is not  
695 sufficient to establish external validity (Dever and Valliant 2014).

696 External validity may also require establishing that the population represented by the VGI is the  
697 same as the population of the full study region. Two examples are provided to illustrate this practical  
698 issue. In both examples, the objective is to estimate the accuracy of a map. For the first example, suppose  
699 that volunteers avoid locations of complex land cover and provide reference data exclusively for locations  
700 that are surrounded by homogeneous land cover. Antoniou et al. (2016) suggest such a strategy may be  
701 beneficial when using photographs to avoid difficulties of determining the ground condition. Because  
702 homogeneous regions are typically more likely to be classified correctly, the accuracy estimates produced  
703 from such data would be expected to have higher accuracy than is true of the study region as a whole.  
704 Consequently external validity of these data would be suspect because the estimates based on the non-  
705 probability sample would not be generalizable to the target population. As a second example, suppose  
706 because of convenient access the VGI data have been collected primarily at locations near roads.  
707 Evaluating external validity would then require determining whether accuracy near roads was equivalent  
708 to accuracy distant from roads.

709 Verifying external validity of VGI may be extremely challenging and in some cases impossible  
710 (Dever and Valliant 2014). Verification requires comparing characteristics of the VGI data with

711 characteristics of the full study region. Consider the example of VGI data concentrated along roads. To  
712 establish that accuracy does not vary with distance from a road, we could collect additional reference  
713 data distant from roads based on a probability sampling design, and compare the accuracy estimates  
714 from this sample to accuracy estimates for sample data constrained to locations near roads. But the  
715 additional effort to obtain the sample data distant from roads would negate much of the value of VGI  
716 for reducing the cost of accuracy assessment. That is, to definitively establish the equivalence of  
717 accuracy near roads to accuracy distant from roads, we may need a large probability sample, and the  
718 primary value of VGI is to reduce the cost and effort of collecting sample data.

719       Alternatively, it may be possible to cite previous studies to establish external validity. For example,  
720 if previous research has demonstrated that distance from a road is not strongly related to accuracy, we  
721 would have some assurance of external validity to support use of VGI data collected preferentially near  
722 roads. In general, to more fully exploit the potential benefit of VGI, it may be necessary to document  
723 typical features of VGI that would commonly need to be addressed to establish external validity and  
724 then conduct the necessary studies to inform the decision of whether external validity is tenable.  
725 Distance from road, characteristics of volunteers, and complexity of landscape are just a few examples  
726 of features that might be explored to determine whether characteristics of populations (e.g., accuracy)  
727 differ by these features. If in general there are no such differences, external validity of non-probability  
728 sample data is supported to some degree. Developing a cohesive strategy to design and conduct such  
729 studies for a broadly applicable assessment of external validity of VGI would likely require a major  
730 research initiative.

731

### 732 **5.3 VGI and Model-Based Inference**

733       Model-based inference is not predicated on probability sampling so it is a potentially attractive  
734 option for using VGI data that did not originate from a probability sampling design. Model-based

735 inference requires specification of a model that relates  $y_u$  to a set of covariates (predictors) available for  
736 the full population (Valliant et al. 2000). Developing appropriate models and evaluating the underlying  
737 assumptions may be difficult and time-consuming (Baker et al. 2013) with the difficulties exacerbated by  
738 the fact that in most surveys, numerous estimates are produced from a single sample. In the case of  
739 VGI, estimates of accuracy and area for several land-cover or land-cover change types will typically be of  
740 interest, and each of these estimates may be desired for several subregions within the target region of  
741 interest. A model will need to be developed and assumptions evaluated for all estimates as a model  
742 that works well for some estimates may not work well for others. An additional challenge to the model-  
743 based approach is that non-probability samples may have an inherent selection bias, so a substantial risk  
744 exists that the distribution of important covariates in the sample will differ from the distribution of these  
745 covariates in the target population (Baker et al. 2013). Methods to account for preferential sampling  
746 (e.g., Diggle et al. 2010) in a model-based framework may be considered in such cases of non-probability  
747 sampling.

748 Numerous model-based methods can be applied to non-probability samples and evaluating the  
749 utility of model-based methods is case specific because it is difficult to ascribe general properties to  
750 these methods (Baker et al. 2013). An advantage of probability sampling and design-based inference is  
751 that a standard general approach is used to produce the complete array of estimates (see Section 2.1).  
752 Yet another challenge of model-based inference and non-probability sampling is how to define and  
753 quantify uncertainty. A widely accepted measure of precision does not exist for estimates from non-  
754 probability samples (Baker et al. 2013, p.97), whereas the standard error (or appropriately scaled  
755 version of standard error) is generally accepted for quantifying precision of estimates in design-based  
756 inference. Clearly, some of the cost savings achieved by non-probability sampling is lost due to the  
757 more complex analyses needed to develop models and test their assumptions (Baker et al. 2013).  
758 Because model-based inference encompasses an array of methods, establishing transparency of the

759 methodology is also more demanding because it is necessary to describe the specific model-based  
760 approach used and the possible limitations of inference uniquely associated with that approach (Baker  
761 et al. 2013, p.100).

762

## 763 **6. Discussion**

764 The increasing availability of large quantities of data obtained via non-probability sampling has  
765 garnered interest of survey methodologists in a variety of subject areas, so it is relevant to examine  
766 issues addressed in the broader survey sampling literature that go beyond just use of VGI in the remote  
767 sensing context. For example, internet surveys comprised of volunteer opt-in panels that use social  
768 media to extract information result in large quantities of data that are obtained quickly and conveniently  
769 but via a selection protocol that has no underlying probability sampling design. Review articles by Baker  
770 et al. (2013) and Elliott and Valliant (2017) provide an excellent general overview of methods and issues  
771 affecting inference when using data from such non-probability samples. In the broad context of survey  
772 sampling, the conventional practice of relying on design-based inference has been questioned because  
773 of the tremendous increase in non-response rates. Even if a probability sampling design is  
774 implemented, severe non-response will make the application of design-based inference questionable  
775 (Baker et al. 2013). Fortunately, in land-cover studies non-response is generally not a major problem.  
776 The availability of remote sensing platforms usually allows us to obtain the necessary observations that  
777 might otherwise be very difficult if a ground visit were required. Non-response rates are typically very  
778 small in accuracy assessment and area estimation applications so the dilemma of severe non-response  
779 that impacts current survey practice in other fields of application is typically not a problem in land-cover  
780 studies.

781 Ensuring accurate observations ( $y_u$ ) is perhaps the most challenging aspect of using VGI because it  
782 depends on the volunteers to provide good quality data. Accurate interpretation of reference labels for

783 land cover or land-cover change is challenging even for trained experts so label quality of VGI data needs  
784 to be scrutinized closely. A great deal of effort has been invested in improving and evaluating the  
785 quality of VGI used in land-cover studies, including the assessment of traditional quality measures such  
786 as positional, thematic or temporal accuracy (Fonte et al. 2017a), the development of new quality  
787 indicators that are applicable specifically to VGI (Meek et al. 2014; Antoniou and Skopeliti 2015;  
788 Senaratne et al. 2017), and even combinations of indicators (Bishr and Mantelas 2008; Jokar Arsanjani et  
789 al. 2015). The investment in these methods will not only yield better quality VGI data but may also  
790 contribute to improved data quality and assessment procedures applicable to reference data obtained  
791 by experts.

792 Baker et al. (2013) make the helpful distinction between “describers” whose purpose is to describe  
793 the population and “modelers” whose purpose is to characterize relationships between variables.  
794 Accuracy assessment and area estimation applications typically fall within the “describer” class because  
795 of the strong focus on descriptive parameters such as user’s and producer’s accuracies of the different  
796 classes and the area or proportion of area of the land-cover or land-cover change classes. Describers  
797 generally rely on probability sampling because of the importance of representing the target population.  
798 Elliott and Valliant (2017, p.262) provide a strong statement in support of probability sampling for  
799 descriptive objectives: “... when critical estimates of descriptive quantities such as means, quantiles or  
800 cell probabilities are required, nonprobability designs should be avoided or utilized only when it is  
801 reasonably certain that there are available covariates in both datasets related to the nonprobability  
802 selection mechanism that can be used to appropriately incorporate information from the nonprobability  
803 sample. If a sufficiently large probability sample is available for estimating descriptive statistics,  
804 methods to incorporate nonprobability data are likely not warranted.”

805 Although design-based inference requires a probability sampling design, it is not reasonable to  
806 assert a recommendation that probability sampling must always be used. Other considerations such as

807 cost and “fit for purpose “may be relevant, the latter including dimensions such as “accuracy, timeliness,  
808 and accessibility” (Baker et al. 2013, p. 98). A quote from Kish (1965, pp. 28-29) extracted by Baker et al.  
809 (2013, p.92) has direct bearing on this issue: “No clear rule exists for deciding exactly when probability  
810 sampling is necessary, and what price should be paid for it ... Probability sampling for randomization is  
811 not a dogma, but a strategy, especially for large numbers.” Probability sampling offers the strong  
812 advantage that it provides the basis for rigorous design-based inference, but there may be exceptional  
813 cases in which fit for purpose criteria will be such that VGI from a non-probability sample will suffice.  
814 While an unmistakable conclusion from our assessment of VGI for use in design-based inference is that  
815 probability sampling should be used, we recognize that occasionally circumstances may exist where not  
816 following this recommendation is justifiable.

817 VGI has great potential value within remote sensing beyond its use to produce accuracy and  
818 area estimates within design-based inference. For example, VGI can greatly augment traditional sources  
819 of training data used in the classification algorithms of land cover and land use maps. The exact design  
820 of the training stage of a supervised classification should, however, be highly classifier-specific as  
821 classifiers vary greatly in how they use the training set. While conventional statistical classifiers may  
822 benefit from the use of a probability sample in the acquisition of training statistics to obtain a  
823 representative and unbiased description of each class, other classifiers, such as machine learning  
824 classifiers, may require only very small and distinctly non-random sample. Thus, for example, an  
825 effective approach to training data acquisition for a classification by a support vector machine may be to  
826 direct citizens to a small number of highly atypical training sites (Pal and Foody 2012). Classifiers also  
827 vary in their sensitivity to mis-labeling of training cases (Foody et al. 2016) which may be relevant if VGI  
828 is to be used.

829 Land cover data from several Geo-Wiki campaigns are now available in the openly accessible  
830 repository Pangaea and these data could be used as training data (Fritz et al. 2017; Laso Bayas et al.

831 2017). VGI is also useful in the development of hybrid land-cover maps, where methods such as  
832 geographically weighted regression can use VGI to determine the most appropriate land cover class at a  
833 given location among several existing products. Such an approach has been demonstrated in the  
834 development of global land cover and forest masks (Schepaschenko et al. 2015; See et al. 2015). Finally,  
835 VGI can provide a preliminary check on the accuracy of a land-cover product and guide the collection of  
836 additional training data in areas where there is visual evidence of confusion between land-cover classes.

837

## 838 **7. Summary**

839 The increasing availability and quantity of VGI has generated great interest in how these data might  
840 be used in applications requiring land-cover data, specifically area estimation and map accuracy  
841 assessment. Scientifically credible use of VGI raises many of the same issues related to inference that  
842 McRoberts (2011) discussed pertaining to use of land-cover maps, stating that “...rules must be  
843 rigorously followed to produce valid scientific inferences.” The requirements for using VGI in rigorous  
844 design-based inference are identifiable from the analysis protocol (Sec. 3.1) used to produce the area  
845 and map accuracy estimates. Specifically, the estimates are derived from totals, and the Horvitz-  
846 Thompson estimator provides an unbiased estimator of a population total if the response design  
847 generates accurate observation of the attribute or measurement of interest ( $y_u$ ) and the sampling design  
848 is such that the inclusion probabilities ( $\pi_u$ ) are known. If  $y_u$  is accurate and  $\pi_u$  is known then we can  
849 produce unbiased estimators of the totals that form the basis for accuracy and area estimates. We  
850 reviewed recent literature describing methods for obtaining VGI and assessing its quality (Sec. 3.2), and  
851 we anticipate that ongoing research will improve reference data quality whether collected by volunteers  
852 within a VGI framework or by expert interpreters.

853 The primary focus of this article has been on the sampling design issues related to using VGI in  
854 design-based inference, with attention addressing three primary cases: 1) VGI data are from a

855 probability sampling design; 2) VGI data from a non-probability sampling design are combined with data  
856 from a probability sampling design; and 3) the only data available are VGI data from a non-probability  
857 sampling design. The most direct approach to ensure that design-based inference can be invoked is to  
858 specify that the VGI data will be collected at locations (sample units) selected by a probability sampling  
859 design (“active VGI”). Implementing a probability sampling design ensures that the inclusion  
860 probabilities ( $\pi_u$ ) for the sampled units are known and thus the corresponding estimation weights  
861 ( $w_u=1/\pi_u$ ) required for the analysis are known. The more common situation is that the VGI data do not  
862 originate from a probability sampling design. Implementing design-based inference in this situation  
863 requires combining the VGI data with data obtained from a probability sampling design, and the benefit  
864 of the VGI data is to reduce the standard errors of the accuracy or area estimates. Two approaches for  
865 combining VGI with a probability sample are to treat the VGI as a certainty stratum (i.e., set  $\pi_u=1$  for  
866 each unit from the VGI sample) or to use the VGI to create an auxiliary variable for the population and  
867 incorporate this variable in a model-assisted estimator. The certainty stratum approach is the more  
868 promising of these two options particularly if a large proportion of the population is covered by VGI. For  
869 land-cover studies the model-assisted estimator use of VGI likely will also incorporate maps produced  
870 from remote sensing imagery.

871 If VGI data collected from a non-probability sampling design are the only data available, rigorous  
872 design-based inference is not available. Estimates of accuracy and area can be produced using the same  
873 estimator formulas of design-based inference by defining pseudo-estimation weights based on treating  
874 the VGI as if a stratified random sample had been implemented. Estimates produced in this fashion  
875 mimic the proportional representation of the feature of the population used to create the pseudo-  
876 weights. However, in contrast to the case where the weights are the inverse of known inclusion  
877 probabilities from a probability sampling design, the estimates based on pseudo-weights require the  
878 additional step of verifying that the condition of external validity is satisfied. External validity requires

879 that the population for which the VGI data are representative must have the same characteristics (e.g.,  
880 model relationships) as the full population that is the target of inference. Establishing external validity is  
881 often impractical so the pseudo-weight approach to using VGI from a non-probability sample will have  
882 limited utility. Model-based inference is perhaps the more promising avenue for using VGI from non-  
883 probability samples. Explication of model-based methods and specific example applications of accuracy  
884 and area estimation (McRoberts 2006; Magnussen 2015) are needed to make model-based inference  
885 more accessible to practitioners.

886 Invoking design-based inference as the scientific basis to support the validity of inference for  
887 estimating area and map accuracy from sample data imposes the requirement that the sampling and  
888 estimation protocols implemented must satisfy certain conditions. As is apparent from the methods and  
889 discussion presented in this article, the requirement of a probability sampling design places fairly strong  
890 restrictions on how VGI can be used in design-based inference. The methods presented in this article for  
891 incorporating VGI in design-based inference expand the potential utility of this growing body of data for  
892 applications of accuracy assessment and area estimation.

893

#### 894 **Acknowledgments**

895 This research was supported in part by the Portuguese Foundation for Science and Technology (FCT)  
896 under project grant UID/MULTI/00308/2013 (CF); EU-funded FP7 project CrowdLand No. 617754 and  
897 the Horizon2020 LandSense project No. 689812 (LS); Cooperative Agreement G12AC20221 provided by  
898 the United States Geological Survey and NASA Carbon Monitoring System program grant  
899 NNX13AP48G (SS). We thank the reviewers for their constructive comments that led to improvements  
900 in the manuscript.

901

#### 902 **References**

903

904 Antoniou, V., Morley, J., and Haklay, M. (2010). Web 2.0 geotagged photos: Assessing the spatial  
905 dimension of the phenomenon. *Geomatica*, 64, 99–110.

906

907 Antoniou, V., Fonte, C. C., See, L., Estima, J., Arsanjani, J. J., Lupia, F., Minghini, M., Foody, G., and Fritz,  
908 S. (2016). Investigating the feasibility of geo-tagged photographs as sources of land cover input data.  
909 *ISPRS International Journal of Geo-Information*, 5, 64.

910

911 Antoniou, V., and Skopeliti, A. (2015). Measures and indicators of VGI quality: An overview. In: *ISPRS*  
912 *Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*. Presented at the  
913 ISPRS Geospatial Week 2015, ISPRS Annals, La Grande Motte, France, pp. 345–351.

914 Baker, R., Brick, J. M., Bates, N. A., Battaglia, M., Couper, M. P., Dever, J. A., Gile, K. J., and Tourangeau,  
915 R. (2013). Summary report of the AAPOR Task Force on non-probability sampling. *Journal of Survey*  
916 *Statistics and Methodology*, 1, 90-105.

917  
918 Bird, T. J., Bates, A. E., Lefcheck, J. S., Hill, N. A., Thomson, R. J., Edgar, G. J., Stuart-Smith, R. D.,  
919 Wotherspoon, S., Krkosek, M., Stuart-Smith, J. F., Pecl, G. T., Barrett, N., and Frusher, S. (2014).  
920 Statistical solutions for error and bias in global citizen science datasets. *Biological Conservation*, 173,  
921 144-154.

922  
923 Bishr, M., and Mantelas, L. (2008). A trust and reputation model for filtering and classifying knowledge  
924 about urban growth. *GeoJournal*, 72, 229–237.

925  
926 Brus, D. J., and de Gruijter, J. J. (2003). A method to combine non-probability sample data with  
927 probability sample data in estimating spatial means of environmental variables. *Environmental*  
928 *Monitoring and Assessment*, 83, 303-317.

929  
930 Card, D. H. (1982). Using known map category marginal frequencies to improve estimates of thematic  
931 map accuracy. *Photogrammetric Engineering & Remote Sensing*, 48, 431-439.

932 Ciepluch, B., Jacob, R., Mooney, P., and Winstanley, A. (2010). Comparison of the accuracy of  
933 OpenStreetMap for Ireland with Google Maps and Bing Maps. *Proceedings of the Ninth International*  
934 *Symposium on Spatial Accuracy Assessment in Natural Resources and Environmental Sciences*, 20-23  
935 July, 2010, 337.

936 Clark, M.L. and Aide, T.M. (2011). Virtual Interpretation of Earth Web-Interface Tool (VIEW-IT) for  
937 collecting land-use/land-cover reference data. *Remote Sensing*, 3, 601–620.

938  
939 Cochran, W. G. (1977). *Sampling Techniques*, 3rd ed. (New York: John Wiley & Sons).

940  
941 De Gruijter, J. J., and Ter Braak, C. J. F. (1990). Model-free estimation from spatial samples: A reappraisal  
942 of classical sampling theory. *Mathematical Geology*, 22, 407-415.

943  
944 Dever, J. A., Rafferty, A., and Valliant, R. (2008). Internet surveys: Can statistical adjustments eliminate  
945 coverage bias? *Survey Research Methods*, 2, 47-62.

946  
947 Dever, J. A., and Valliant, R. (2014). Estimation with non-probability surveys and the question of external  
948 validity. *Proceedings of Statistics Canada Symposium 2014*, 8 pp.

949

950 Deville, J. C., and Särndal, C. E. (1992). Calibration estimators in survey sampling. *Journal of the*  
951 *American Statistical Association*, 87, 376-382.

952

953 Diggle, P., Menezes, R., and Su, T. (2010). Geostatistical inference under preferential sampling. *Applied*  
954 *Statistics*, 59, 191-232.

955

956 Diggle, P. J., Tawn, J. A., and Moyeed, R. A. (1998). Model-based geostatistics. *Applied Statistics*, 47, 299-  
957 350.

958

959 Elliott, M. R., and Valliant, R. (2017). Inference for nonprobability samples. *Statistical Science*, 32, 249-  
960 264.

961

962 Elwood, S., Goodchild, M.F., and Sui, D. (2013). Prospects for VGI research and the emerging fourth  
963 paradigm. *Crowdsourcing geographic knowledge: Volunteered geographic information (VGI) in theory*  
964 *and practice*, Springer, Netherlands, pp. 361-375.

965

966 Fonte, C. C., Bastin, L., See, L., Foody, G., and Lupia, F. (2015). Usability of VGI for validation of land cover  
967 maps. *International Journal of Geographical Information Science*, 29, 1269-1291.

968

969 Fonte, C.C., Antoniou, V., Bastin, L., Estima, J., Arsanjani, J.J., Laso-Bayas, J.-C., See, L., and Vatseva, R.  
970 (2017a). Assessing VGI data quality, in: Foody, G.M., See, L., Fritz, S., Fonte, C.C., Mooney, P., Olteanu-  
971 Raimond, A.-M., Antoniou, V. (Eds.), *Mapping and the Citizen Sensor*. Ubiquity Press, London, UK, pp.  
972 137-164.

973

974 Fonte, C. C., Patriarca, J. A., Minghini, M., Antoniou, V., See, L., and Brovelli, M. A. (2017b). Using  
975 OpenStreetMap to create land use and land cover maps: Development of an application. In: *Volunteered*  
976 *Geographic Information and the Future of Geospatial Data* (Ed.: Campelo, C. E. C., Bertolotto, M., and  
977 Corcoran, P.). IGI Global. DOI: 10.4018/978-1-5225-2446-5.ch007

977

978 Foody, G. M. (2009). The impact of imperfect ground reference data on the accuracy of land cover  
979 change estimation. *International Journal of Remote Sensing*, 30, 3275-3281.

979

980 Foody, G.M. (2010). Assessing the accuracy of land cover change with imperfect ground reference data.  
981 *Remote Sensing of Environment*, 114, 2271-2285.

981

982 Foody, G.M., Pal, M., Rocchini, D., Garzon-Lopez, C.X. and Bastin, L. (2016). The sensitivity of mapping  
983 methods to reference data quality: Training supervised image classifications with imperfect reference  
984 data. *ISPRS International Journal of Geo-Information*, 5(11), p.199.

985

986 Foody, G. M., See, L., Fritz, S., Van der Velde, M., Perger, C., Schill, C., and Boyd, D. S. (2013). Assessing  
987 the accuracy of volunteered geographic information arising from multiple contributors to an internet  
988 based collaborative project. *Transactions in GIS*, 17, 847–860.

989 Foody, G.M., See, L., Fritz, S., Velde, M. van der, Perger, C., Schill, C., Boyd, D.S., and Comber, A. (2015).  
990 Accurate attribute mapping from Volunteered Geographic Information: Issues of volunteer quantity and  
991 quality. *The Cartographic Journal*, 52, 336–344.

992 Fritz, S., McCallum, I., Schill, C., Perger, C., See, L., Schepaschenko, D., van der Velde, M., Kraxner, F., and  
993 Obersteiner, M. (2012). Geo-Wiki: An online platform for improving global land cover. *Environmental*  
994 *Modelling & Software*, 31, 110–123.

995 Fritz, S., See, L., van der Velde, M., Nalepa, R.A., Perger, C., Schill, C., McCallum, I., Schepaschenko, D.,  
996 Kraxner, F., Cai, X., Zhang, X., Ortner, S., Hazarika, R., Cipriani, A., Di Bella, C., Rabia, A.H., Garcia, A.,  
997 Vakolyuk, M., Singha, K., Beget, M.E., Erasmi, S., Albrecht, F., Shaw, B., and Obersteiner, M. (2013).  
998 Downgrading recent estimates of land available for biofuel production. *Environmental Science and*  
999 *Technology*, 47, 1688–1694.

1000 Fritz, S., See, L., Perger, C., McCallum, I., Schill, C., Schepaschenko, D., Duerauer, M., Karner, M., Dresel,  
1001 C., Laso Bayas, J.C., Lesiv, M., Moorthy, I., Salk, C., Danylo, O., Sturn, T., Albrecht, F., You, L., Kraxner, F.  
1002 and Obersteiner, M. (2017). A global dataset of crowdsourced land cover and land use reference data.  
1003 *Scientific Data*, 4:170075. doi: 10.1038/sdata.2017.75.

1004 Gallego, F. J. (2004). Remote sensing and land cover area estimation. *International Journal of Remote*  
1005 *Sensing*, 25, 3019-3047.

1006 Girres, J.-F., and Touya, G. (2010). Quality assessment of the French OpenStreetMap dataset.  
1007 *Transactions in GIS*, 14, 435–459.

1008 Goodchild, M. F. (2007). Citizens as sensors: The world of volunteered geography. *GeoJournal*, 69, 211–  
1009 221.

1010

1011 Gregoire, T. G. (1998). Design-based and model-based inference in survey sampling: appreciating the  
1012 difference. *Canadian Journal of Forest Research*, 28, 1429-1447.

1013

1014 Haklay, M. (2010). How good is volunteered geographical information? A comparative study of  
1015 OpenStreetMap and Ordnance Survey datasets. *Environment and Planning B: Planning and Design*,  
1016 37,682–703.

1017

1018 Horvitz, D. G., and Thompson, D. J. (1952). A generalization of sampling without replacement from a  
1019 finite universe. *Journal of the American Statistical Association*, 47, 663-685.

1020

1021 Hou, D., Chen, J., Wu, H., Li, S., Chen, F., and Zhang, W. (2015). Active collection of land cover sample  
1022 data from geo-tagged web texts. *Remote Sensing*, 7, 5805-5827.

1023

1024 Isaaks, E. H., and Srivastava, R. M. (1989). *An Introduction to Applied Geostatistics*, Oxford University  
1025 Press, New York.

1026

- 1027 Iwao, K., Nishida, K., Kinoshita, T., and Yamagata, Y. (2006). Validating land cover maps with Degree  
 1028 Confluence Project information. *Geophysical Research Letters*, 33, L23404, doi:10.1029/2006GL027768.  
 1029
- 1030 Jokar Arsanjani, J. and Bakillah, M. (2015). Understanding the potential relationship between the socio-  
 1031 economic variables and contributions to OpenStreetMap. *International Journal of Digital Earth*, 8, 861–  
 1032 876.
- 1033 Kish, L. (1965). *Survey Sampling*. Wiley, New York.
- 1034 Laso Bayas, J.C., Lesiv, M., Waldner, F., Schucknecht, A., Duerauer, M., See, L. et al. (2017). A global  
 1035 reference database of crowdsourced cropland data collected using the Geo-Wiki platform. *Scientific*  
 1036 *Data*. 4: 170136. doi:10.1038/sdata.2017.136.
- 1037 Laso Bayas, J.-C., See, L., Fritz, S., Sturn, T., Perger, C., Duerauer, M., Karner, M., Moorthy, I.,  
 1038 Schepaschenko, D., Domian, D., and McCallum, I. (2017). Crowdsourcing in-situ data on land cover and  
 1039 land use using gamification and mobile technology. *Remote Sensing* 8(11), 905. doi:10.3390/rs8110905.
- 1040 Loosveldt, G., and Sonck, N. (2008). An evaluation of the weighting procedures for an online access  
 1041 panel survey. *Survey Research Methods*, 2, 93-105.  
 1042
- 1043 Magnussen, S. (2015). Arguments for a model-dependent inference? *Forestry*, 88, 317-325.
- 1044 Martino, L., Palmieri, A., and Gallego, J. (2009). Use of auxiliary information in the sampling strategy of a  
 1045 European area frame agro-environmental survey. Available from:  
 1046 [http://ec.europa.eu/eurostat/documents/205002/769457/LUCAS2009\\_S2-Sampling\\_20090000.pdf](http://ec.europa.eu/eurostat/documents/205002/769457/LUCAS2009_S2-Sampling_20090000.pdf).  
 1047 Eurostat.
- 1048 McRoberts, R. E. (2006). A model-based approach to estimating forest area. *Remote Sensing of*  
 1049 *Environment*, 103, 56-66.
- 1050 McRoberts, R. E. (2010). Probability- and model-based approaches to inference for proportion forest  
 1051 using satellite imagery as ancillary data. *Remote Sensing of Environment*, 114, 1017-1025.
- 1052 McRoberts, R. E. (2011). Satellite image-based maps: Scientific inference or pretty pictures? *Remote*  
 1053 *Sensing of Environment*, 115, 715–724.
- 1054 Meek, S., Jackson, M.J., and Leibovici, D.G. (2014). A flexible framework for assessing the quality of  
 1055 crowdsourced data. In: Huerta, J., Schade, S., and Granell, C. (Eds.), *Connecting a Digital Europe through*  
 1056 *Location and Place: Proceedings of the AGILE'2014 International Conference on Geographic Information*  
 1057 *Science*. Presented at the AGILE'2014 International Conference on Geographic Information Science,  
 1058 AGILE Digital Editions, Castellón, Spain.
- 1059 Meek, S., Jackson, M., and Leibovici, D.G. (2016). A BPMN solution for chaining OGC services to quality  
 1060 assure location-based crowdsourced data. *Computers & Geosciences*, 87, 76–83.

- 1061 Mullen, W. F., Jackson, S. P., Croitoru, A., Crooks, A., Stefanidis, A., and Agouris, P. (2014). Assessing the  
1062 impact of demographic characteristics on spatial error in volunteered geographic information features.  
1063 *GeoJournal*, 80, 587-605.
- 1064 Olofsson, P., Foody, G. M., Stehman, S. V., and Woodcock, C. E. (2013). Making better use of accuracy  
1065 data in land change studies: Estimating accuracy and area and quantifying uncertainty using stratified  
1066 estimation. *Remote Sensing of Environment*, 129, 122-131.
- 1067  
1068 Olofsson, P., Foody, G. M., Herold, M., Stehman, S. V., Woodcock, C. E., and Wulder, M. A. (2014). Good  
1069 practices for estimating area and assessing accuracy of land change. *Remote Sensing of Environment*,  
1070 148, 42-57.
- 1071  
1072 Overton, J. M., Young, T. C., and Overton, W. S. (1993). Using 'found' data to augment a probability  
1073 sample: Procedure and case study. *Environmental Monitoring and Assessment*, 26, 65-83.
- 1074  
1075 Pal, M. and Foody, G.M. (2012). Evaluation of SVM, RVM and SMLR for accurate image classification with  
1076 limited ground data. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*,  
1077 5, 1344-1355.
- 1078 Sannier, C., McRoberts, R. E., Fichet, L.-V., & Makaga, E. M. K. (2014). Using the regression estimator  
1079 with Landsat data to estimate proportion forest cover and net proportion deforestation in Gabon.  
1080 *Remote Sensing of Environment*, 151, 138-148.
- 1081 Särndal, C. E., Swensson, B., and Wretman, J. (1992). *Model-Assisted Survey Sampling*. Springer-Verlag,  
1082 New York.
- 1083 Schepaschenko, D., See, L., Lesiv, M., McCallum, I., Fritz, S., Salk, C., Moltchanova, E., Perger, C.,  
1084 Shchepashchenko, M., Shvidenko, A., Kovalevskiy, S., Gilitukha, D., Albrecht, F., Kraxner, F., Bun, A.,  
1085 Maksyutov, S., Sokolov, A., Dürauer, M., Obersteiner, M., Karminov, V., and Ontikov, P. (2015).  
1086 Development of a global hybrid forest mask through the synergy of remote sensing, crowdsourcing and  
1087 FAO statistics. *Remote Sensing of Environment*, 162, 208–220.
- 1088 See, L., Comber, A., Salk, C., Fritz, S., van der Velde, M., Perger, C., Schill, C., McCallum, I., Kraxner, F.,  
1089 and Obersteiner, M. (2013). Comparing the quality of crowdsourced data contributed by expert and  
1090 non-experts. *PLoS ONE* 8, e69958. doi:10.1371/journal.pone.0069958
- 1091 See, L., Fritz, S., Perger, C., Schill, C., McCallum, I., Schepaschenko, D., Dürauer, M., Sturn, T., Karner,  
1092 M., Kraxner, F., and Obersteiner, M. (2015). Harnessing the power of volunteers, the internet and  
1093 Google Earth to collect and validate global spatial information using Geo-Wiki. *Technological Forecasting  
1094 and Social Change* 98, 324–335.
- 1095 See, L., Mooney, P., Foody, G., Bastin, L., Comber, A., Estima, J., Fritz, S., Kerle, N., Jiang, B., Laakso, M.,  
1096 Liu, H.-Y., Milčinski, G., Nikšič, M., Painho, M., Pósdör, A., Olteanu-Raimond, A.-M., and Rutzinger, M.  
1097 (2016). Crowdsourcing, Citizen Science or Volunteered Geographic Information? The current state of

1098 crowdsourced geographic information. *ISPRS International Journal of Geo-Information*, 5(5), 55. doi:  
1099 10.3390/ijgi5050055

1100 See, L., Laso Bayas, J.C., Schepaschenko, D., Perger, C., Dresel, C., Maus, V., Salk, C., Weichselbaum, J.  
1101 Lesiv, M., McCallum, I., Moorthy, I. and Fritz, S. (2017). LACO-Wiki: A new online land cover validation  
1102 tool demonstrated using GlobeLand30 for Kenya. *Remote Sensing* 9(7), 754; doi:[10.3390/rs9070754](https://doi.org/10.3390/rs9070754).

1103 Senaratne, H., Mobasher, A., Ali, A.L., Capineri, C., and Haklay, M. (2017). A review of volunteered  
1104 geographic information quality assessment methods. *International Journal of Geographical Information  
1105 Science*, 31, 139-167.

1106 Stehman, S. V. (1999). Basic probability sampling designs for thematic map accuracy assessment.  
1107 *International Journal of Remote Sensing*, 20, 2423-2441.  
1108

1109 Stehman, S. V. (2000). Practical implications of design-based sampling inference for thematic map  
1110 accuracy assessment. *Remote Sensing of Environment*, 72, 35-45.  
1111

1112 Stehman, S. V. (2009). Model-assisted estimation as a unifying framework for estimating the area of land  
1113 cover and land-cover change from remote sensing. *Remote Sensing of Environment*, 113, 2455-2462.

1114 Stehman, S. V. (2013). Estimating area from an accuracy assessment error matrix. *Remote Sensing of  
1115 Environment*, 132, 202-211.  
1116

1117 Stehman, S. V., and Czaplewski, R. L. (1998). Design and analysis for thematic map accuracy assessment:  
1118 Fundamental principles. *Remote Sensing of Environment*, 64, 331-344.  
1119

1120 Strahler, A. H., Boschetti, L., Foody, G. M., Friedl, M. A., Hansen, M. C., Herold, M., Mayaux, P.,  
1121 Morisette, J. T., Stehman, S. V., and Woodcock, C. E. (2006). Global land cover validation:  
1122 Recommendations for evaluation and accuracy assessment of global land cover maps, EUR 22156 EN –  
1123 DG, Office for Official Publications of the European Communities, Luxembourg, 48 pp.

1124 Tsendbazar, N.E., de Bruin, S., and Herold, M. (2015). Assessing global land cover reference datasets for  
1125 different user communities. *ISPRS Journal of Photogrammetry and Remote Sensing*, 103, 93–114.

1126 Valliant, R. (2013). Comment on “Summary report of the AAPOR Task Force on non-probability  
1127 sampling.” *Journal of Survey Statistics and Methodology*, 1, 105-111.

1128 Valliant, R., and Dever, J. A. (2011). Estimating propensity adjustments for volunteer web surveys.  
1129 *Sociological Methods & Research*, 40, 105-137.

1130 Valliant, R., Dorfman, A. H., and Royall, R. M. (2000). *Finite Population Sampling and Inference: A  
1131 Prediction Approach*, John Wiley & Sons, Inc., New York.  
1132

1133

1134 LIST OF FIGURE CAPTIONS

1135 Figure 1.

1136 Figure 2.

1137 Figure 3.

1138