

Combining imperfect automated annotations of  
underwater images with human annotations to obtain  
precise and unbiased population estimates

Jui-Han Chang<sup>1</sup>, Deborah R. Hart<sup>1\*</sup>, Burton V. Shank<sup>1</sup>, Scott M. Gallagher<sup>2</sup>,  
Peter Honig<sup>2</sup>, and Amber D. York<sup>2</sup>

<sup>1</sup>Northeast Fisheries Science Center, National Marine Fisheries Service,  
166 Water St., Woods Hole, MA 02543, USA

<sup>2</sup>Biology Department, Woods Hole Oceanographic Institution,  
Woods Hole, MA 02543, USA

\*Email: [deborah.hart@noaa.gov](mailto:deborah.hart@noaa.gov)

August 10, 2016

## 2 **Abstract**

3 Optical methods for surveying populations are becoming increasingly popular. These meth-  
4 ods often produce hundreds of thousands to millions of images, making it impractical to an-  
5 alyze all the images manually by human annotators. Computer vision software can rapidly  
6 annotate these images, but their error rates are often substantial, vary spatially and are  
7 autocorrelated. Hence, population estimates based on the raw computer automated counts  
8 can be seriously biased. We evaluated four estimators that combine automated annotations  
9 of all the images with manual annotations from a random sample to obtain (approximately)  
10 unbiased population estimates, namely: ratio, offset, and linear regression estimators as well  
11 as the mean of the manual annotations only. Each of these estimators was applied either  
12 globally or locally (i.e., either all data were used or only those near the point in question, to  
13 take into account spatial variability and autocorrelation in error rates). We also investigated  
14 a simple stratification scheme that splits the images into two strata, based on whether the  
15 automated annotator detected no targets or at least one target. The 16 methods result-  
16 ing from a combination of four estimators, global or local estimation, and one stratum or  
17 two strata, were evaluated using simulations and field data. Our results indicated that the  
18 probability of a false negative is the key factor determining the best method, regardless of  
19 the probability of false positives. Stratification was the most effective method in improving  
20 the accuracy and precision of the estimates, provided the false negative rate was not too  
21 high. If the probability of false negatives are low, stratified estimation with the local ratio  
22 estimator or local regression (essentially geographically weighted regression) are best. If the  
23 probability of false negatives are high, no stratification with a simple global linear regression  
24 or simply the manual sample mean alone is recommended.

25

26 **Keywords:** Underwater imagery; Computer vision; Population estimation; Scallop; Geo-  
27 graphically weighted regression

# 1 Introduction

Underwater optical surveys of fish and invertebrate populations are becoming increasingly common (e.g., Davis et al., 1992; Gallager et al., 2005; Howland et al., 2006; Yoklavich et al., 2007; Rosenkranz et al., 2008; Taylor et al., 2008; Tolimieri et al., 2008; Singh et al., 2013; Gallager et al., 2014). Such surveys have numerous advantages over traditional surveys using fishing gear, including being able to observe populations at all scales under natural conditions, and detection efficiency that potentially approaches 100%.

Optical surveys often generate hundreds of thousands to millions of images. Manually annotating all of the images (i.e., having people identifying the targets of interest in each image) would thus often be impractical. The traditional statistical approach to this problem would be to only manually annotate a sample of the images and obtain inferences on the population (which for our purposes is defined as the targets contained in all of the collected images) based on the sample. Alternatively, computer vision software can produce “automated annotations” that identify the targets in every image. However, automated annotators can make errors, both because they may not detect some targets (“false negatives”) and because the annotator mistakenly identifies some objects (“distractors”) as targets when they are not (“false positives”). Thus, analyses based on the raw automated counts can be seriously biased. Errors from automated annotations are often autocorrelated and spatially non-stationary due to, for example, a certain region having high densities of distractors or reduced visibility. Manual annotations of a sample of the images can help detect and correct for errors by the automated annotators, in which case the goal is to produce estimators for the population, based on the combination of automated and manual annotations that are more efficient than using the manual annotations alone (i.e., the variances of estimators are less than the variance of the sample mean of the manual images), as well as being at least approximately unbiased.

Although there have been numerous studies devoted to automated detection and classification of marine organisms (e.g., Culverhouse et al., 2006; Marcos et al., 2008; Spampinato et al., 2008; Beijbom et al., 2012), these studies usually conclude with estimating confusion matrices or error rates. The final task of obtaining estimates of the population of targets in

57 all images from automated annotations that contain errors has received less attention. Solow  
58 et al. (2001) considered the situation where classification of plankton samples may be in  
59 error, which were corrected by inverting the confusion matrix (see also Hu and Davis, 2006;  
60 Verikas et al., 2015). The problem they considered is simpler than the one we are considering  
61 here because they were only concerned with classification of an object but not its detection,  
62 and because errors were assumed to be stationary and not autocorrelated. Beijbom (2014)  
63 analyzed what we have termed the offset estimator to bias-correct automated counts using  
64 a random sample of manual annotations from a cost reduction point of view.

65 The purpose of this paper is to explore and compare performance of several methods  
66 for estimating population abundance (or biomass) based on automated annotations of all  
67 images combined with manual annotations of a random sample of the images. This study is  
68 motivated by surveys of sea scallops (*Placopecten magellanicus*) using the HabCam (Habitat  
69 Mapping Camera System) towed underwater camera system (Howland et al. 2006; Taylor et  
70 al., 2008; NEFSC, 2014; see Figure 1 for an example of HabCam images of sea scallops and  
71 sand dollars, a common distractor). Computer vision software for detecting sea scallops is  
72 continuing to be developed (Dawkins et al., 2013; Kannappan et al., 2014; Gallager et al., un-  
73 published). The U.S. sea scallop fishery has annual ex-vessel revenue averaging around \$500  
74 million in recent years, so obtaining accurate and precise estimates of sea scallop abundance  
75 is of immediate practical significance.

## 76 **2 Methods, Theory, and Calculation**

### 77 **2.1 Global Population Estimators**

78 We tested four different estimators of population size (i.e., the number of true targets in an  
79 image set) based on a combination of manual and automated annotations. In the following,  
80 it is assumed that each image has been annotated by software, but only a random sample  
81 of  $n$  images out of a total of  $N$  images have been annotated manually, and the manual  
82 annotations are without error (it is straightforward to extend the theory to cases where only  
83 a sample has been annotated by software). Let  $X_i$  and  $Y_i$  be the number of targets detected

84 in the  $i$ th image by the automated and manual annotators, respectively.

Four global estimators for the total number of targets in the images,  $Z$ , are:

Manual sample only:  $Z_m = \bar{Y}N$  (1)

Ratio estimator:  $Z_r = \mu_X N \frac{\bar{Y}}{\bar{X}}$  (2)

Offset estimator:  $Z_o = \sum_{i=1}^N X_i - \frac{N}{n} \sum_{j=1}^n (X_j - Y_j)$  (3)

Regression estimator:  $Z_g = \sum_{i=1}^N \alpha + \beta X_i$  (4)

85 where  $\bar{X}$  and  $\bar{Y}$  are the mean number of targets detected by automated and human annota-  
86 tors in the sample of images that have been manually annotated,  $\mu_X$  is the mean number of  
87 targets over all images detected by the automated annotator, and  $\alpha$  and  $\beta$  in equation (4)  
88 are the intercept and slope obtained by regressing the automated vs. manual annotations.  
89 The last three methods can be considered as ways to adjust, or bias correct, the automated  
90 counts based on the comparison between the automated and manual counts in the sample.  
91 The ratio estimator adjusts the automated counts by a multiplicative constant, the offset es-  
92 timator adjustment by an additive constant, and the regression estimator combines additive  
93 (intercept) and multiplicative (slope) adjustments.

94 Although the ratio estimator (2) is biased, this bias is negligible for all the simulated  
95 datasets because the coefficients of variation of  $\bar{X}$  and  $\bar{Y}$  are both smaller than 0.1 (Cochran,  
96 1977), which should typically be the case because the sample sizes for both the automated  
97 and manual annotations will usually be large. An approximate bias correction can be applied  
98 if this is a concern. The Appendix derives analytically the conditions when the variance of the  
99 ratio estimator applied to a random sample is lower than manual sampling alone. Beijbom  
100 (2014) similarly gave analytic derivations of properties of the offset estimator of a random  
101 sample.

## 102 **2.2 Local Population Estimators**

103 The automated annotator error rate may vary spatially, depending on factors such as water  
104 clarity, substrate type, and the densities of targets and distractors. All these factors, and

105 therefore the automated annotator error rates, are typically spatially autocorrelated. If  
 106 this is the case, it may be more efficient to bias-correct the automated annotations locally,  
 107 rather than using a single global correction as in equations (1)-(4). In addition, the spatial  
 108 distribution of the population is often of interest. If the error rates vary spatially, the  
 109 correction for these errors also needs to vary accordingly to accurately reflect the actual  
 110 distribution of the population.

111 For the local estimators, the correction factor is calculated for each data point, and the  
 112 estimators are similar to the global estimators described above except that only data less  
 113 than a distance, or “bandwidth”,  $h_j$  from the point  $j$  are used, and the data are weighted as  
 114 a decreasing function of the distance from the target data point, using an adaptive bisquare  
 115 distance decay kernel function:

$$w_{(j,k)} = \begin{cases} \left[1 - \left(\frac{d_{(j,k)}}{h_j}\right)^2\right]^2 & d_{(j,k)} \leq h_j \\ 0 & d_{(j,k)} > h_j, \end{cases} \quad (5)$$

116 where  $w_{(j,k)}$  is the weighting factor of point  $k$  that is used to calculate the bias correction  
 117 factor for point  $j$ , and  $d_{(j,k)}$  is the distance between points  $j$  and  $k$ . The bandwidth is  
 118 adapted to the density of the data; it is larger when data are sparser and smaller when  
 119 data are denser. Even though the bandwidth may vary by location, the number of data  
 120 points within the bandwidth is the same across locations. The bandwidth (or number of  
 121 data points to be included at each location) is determined by minimizing the leave-one-out  
 122 cross-validation squared error:

$$CV = \sum_{j=1}^n \left[ Y_j - \hat{Y}_{\neq j}(h_j) \right]^2, \quad (6)$$

123 where  $\hat{Y}_{\neq j}(h_j)$  is the fitted value of  $Y_j$  with the data points where point  $j$  is omitted from  
 124 the estimation process (Guo et al., 2008).

125 The local method for the regression estimator is essentially a form of geographically  
 126 weighted regression (GWR) that is used specifically for situations when the relationship be-  
 127 tween variables differs across space (i.e., spatial non-stationarity and spatial autocorrelation;  
 128 Brunson et al., 2008). Compared to standard (global) regression models where a single pa-

129 parameter set is estimated for the entire dataset, GWR estimates regression parameters that  
130 vary for each data point based on data that is in the local neighborhood of that point.

## 131 2.3 Stratification

132 Population densities from underwater images are often “zero-inflated”, i. e., a high proportion  
133 of photos contain no targets. In such a case, the images can be separated into two strata: one  
134 where no targets were detected by the automated annotator, and the other where at least  
135 one target is detected. Manual annotations are then allocated among the two strata based on  
136 the automated annotations and their overall false negative rates, using approximate Neyman  
137 optimal allocations. For this purpose, the standard deviation of the true target counts in  
138 the zero stratum,  $s_0$ , is:  $\sqrt{Z_0 P_S (1 - P_S)}$ , where  $Z_0$  is the number of targets in the zero  
139 stratum (i.e. the number of false negatives),  $P_S$  is the probability of detecting a target by  
140 the automated annotator, and  $1 - P_S$  is the probability of a false negative. In the simulation,  
141  $Z_0$  and  $P_S$  are known, but in practice, they would have to be estimated either from previous  
142 data or by obtaining a small sample of manual annotations prior to the allocation. The  
143 standard deviation of targets in the non-zero stratum,  $s_1$ , is approximated by the standard  
144 deviation of the automated counts in this stratum. The Neyman optimal allocation is then:

$$n_m = \frac{n N_m s_m}{\sum_{m=0}^1 N_m s_m}, \quad (7)$$

145 where  $n$  is total number of manual sample size, and  $N_m$  is the total number of images in  
146 stratum  $m$ .

## 147 2.4 Simulation Design

148 We tested the performance of the above methods using simulated data. The simulation  
149 design is based on the US sea scallop population characteristics as observed by the HabCam  
150 survey. The simulation domain is 70 km (longitude) by 140 km (latitude), with a 50 m  
151 grid size, roughly corresponding to the density of annotated images in actual data sets. The  
152 spatial distribution of sea scallops is non-stationary due to the influences of physical and  
153 biological environment including current, depth, and predator distributions (Brand, 2006).

154 Therefore, we assumed that the simulated scallop population has large-scale smooth trends  
 155 in its expected mean (first-order effect) that are added to a stationary autocorrelated random  
 156 field (second-order effect; Cressie, 1993). We simulated the variations of global mean density  
 157 using a double logistic function that is constant with latitude but varies with longitude:

$$p(l) = \begin{cases} \frac{1}{1 + \exp(-a(l - b))} & l \leq \frac{l_{max}}{2} \\ \frac{1}{1 + \exp(a(l - b - \frac{l_{max}}{2}))} & l > \frac{l_{max}}{2}, \end{cases} \quad (8)$$

158 where  $l$  is longitude,  $l_{max}$  is the maximum longitude in the surveyed area, and  $a$  and  $b$  are  
 159 the parameters that determine the shape of the logistic curve. The simulated first-order  
 160 effects are high in the middle and decrease logistically toward the left and right edge of the  
 161 simulation domain, which is typical of actual scallop distribution patterns (Hart, 2006). The  
 162 second-order effects were simulated using stationary Gaussian random fields with a spherical  
 163 isotropic covariance structure (Cressie, 1993):

$$\gamma(d) = \begin{cases} 0 & d = 0 \\ c_0 + c_1 \left\{ \frac{3d}{2r} - \frac{1}{2} \left( \frac{d}{r} \right)^3 \right\} & 0 < d \leq r, \\ c_0 + c_1 & d \geq r \end{cases} \quad (9)$$

164 where  $c_0$ ,  $c_1$ , and  $r$  are the nugget, partial sill, and range parameter, respectively. The  
 165 nugget/sill ( $n/s$ ) ratio ( $\frac{c_0}{c_0 + c_1}$ ) determines randomness and  $r$  determines the aggregation  
 166 size of the second-order effects. We chose the simulation parameter values based on estimates  
 167 from the actual HabCam data.

168 To reflect the highly zero-inflated nature of scallop distributions, those locations where the  
 169 sum of the first-order and second-order effects values were smaller than its 90th percentile  
 170 were set to zero. The simulated scallops count for the remaining 10% is simply the sum  
 171 of the first- and second-order effects (Figure 2). The resultant simulated data is patchy,  
 172 zero-inflated, and has a large scale trend along one direction, consistent with actual scallop  
 173 populations. The shape and direction of tracks used to survey the simulated population  
 174 was designed to mimic the actual HabCam survey design, where more effort was put in the



175 middle high density area (Figure 2; NEFSC, 2014). A total of 9,001 photos were simulated  
176 along the track (Figure 2).

177 False positives were simulated by using distractors. The two most common distractors  
178 for sea scallops are sand dollars (*Echinarachnius parma*; Figure 1) and dead scallop shells  
179 (Dawkins et al., 2013; Kannappan et al., 2014). The distribution of sand dollars are typically  
180 independent or negatively correlated with scallops, whereas dead scallop shells would be  
181 expected to be positively related to (live) scallops. The spatial distribution of distractors  
182 were simulated similar to scallops, but the distractor’s patches were assumed larger (larger  
183 range) and less noisy (smaller  $n/s$  ratio) than the scallop target distribution, based on actual  
184 observations of sand dollars (Figure 2).

185 Water visibility may affect automated annotation accuracy by reducing the probability  
186 of detecting a target or a distractor. We simulated water visibility to be trendless but with  
187 spatial autocorrelation. In other words, it is a random field with no first-order effect. It was  
188 assumed to have the same noise level but larger patch size as the distractor (larger range;  
189 Figure 2).

## 190 2.5 Simulation of Automated Count Data

191 The simulated manually annotated data are assumed to have no errors. For the computer  
192 automated counts, each simulated target ( $S$ ) and distractor ( $D$ ) has a probability of being  
193 detected as a target by the automated annotator:

$$P_S = (1 - F1_S)(1 - F2_S) \text{ and } P_D = 1 - (1 - F1_D)(1 - F2_D), \quad (10)$$

194 where the  $F1_S$  and  $F1_D$  are the probabilities of a false negative and false positive with good  
195 water visibility, and  $F2_S$  and  $F2_D$  are the reduced probabilities of detecting targets and  
196 distractors due to water visibility. In our simulations, it is assumed that  $F2_S = F2_D$ . The  
197 simulated total number of targets reported by the automated annotator in the  $i$ th image is:

$$X_i = \sum_{m=1}^M (S_{im} + D_{im}), \quad (11)$$

198 where  $M$  is the total number of objects simulated within image  $i$ ,  $S_{im}$  is the number of

199 correctly identified targets (true positives minus false negatives), and  $D_{im}$  is the number of  
200 distractors incorrectly identified as targets (false positives).

## 201 2.6 Scenarios Tested

202 To understand whether the estimation methods are robust to changes in the environment,  
203 species distributions and the capabilities of the automated annotator, we tested the perfor-  
204 mance of these methods by varying the following quantities:

- 205 (1) Automated annotator’s performance: probability of a false negative/positive ( $F1_S$  and  
206  $F1_D$ ) from 0 to 1 by 0.05;
- 207 (2) Water visibility: good, moderate, or poor (expected value of  $F2 = 0, 0.05, 0.1$ );
- 208 (3) Correlation between scallop and distractor distribution: negative, zero, or positive;
- 209 (4) Degree of spatial autocorrelation of distractors: low, medium, and high;
- 210 (5) Percent of total sample size that was annotated manually: 1%, 3%, 7%, 11%, and 15%.

211 A base case was selected where the water visibility is good, the correlation between the  
212 spatial distribution of scallops and distractors is negative, the spatial autocorrelation of  
213 distractors is medium, and manual annotations were performed on 7% of the photographs.  
214 The base case was then varied for each of the attributes (2)-(5) individually, keeping the  
215 other three at their base case values. Thus, a total of 14 scenarios were simulated. For each  
216 choice of (2)-(5),  $F1_S$  and  $F1_D$  were varied from 0 to 1 by 0.05 increments, as specified in  
217 (1).

218 For all scenarios, scallops have high densities in middle longitudes of the simulation do-  
219 main (simulated using equation 8), and water visibility has no first-order effects. Distractors  
220 have high first-order effects on the left (which used only the second part of the equation 8 on  
221  $l \leq \frac{l_{max}}{2}$  part of the simulation domain), except for the scenarios of zero and positive corre-  
222 lations between scallop and distractor distribution where there are no effects or high effects  
223 in the middle, respectively. The partial sill,  $n/s$  ratio, and range parameter used to simulate  
224 second-order effects are 0.18, 0.6, and 200 for scallops, 0.18, 0.6, and 400 for distractors,  
225 and 0.18, 0.6, and 600 for water visibility. For the scenarios where distractors have high  
226 and low autocorrelation, the  $n/s$  ratio is 0.3 and 0.9, respectively. For the scenarios where

227 water visibility is moderate or poor, the effects of water visibility on the probability of a false  
228 negative and false positive is one or two times, respectively, compared to the corresponding  
229 scenarios of good water visibility.

230 For each scenario, the manual annotation subset was resampled 30 times. For each  
231 iteration, we tested the combinations of the four estimators applied either globally or locally,  
232 and using two strata or one stratum (unstratified) to allocate manual annotations, resulting  
233 in 16 different estimation methods.

234 For stratified estimation, the ratio estimator is undefined in the zero stratum, so the mean  
235 of the manual annotations in this stratum was used instead. Since the offset and regression  
236 estimators reduce to simply taking the mean of the manual annotations in the zero stratum,  
237 all four methods produce the same estimate in this stratum, so any differences among the  
238 methods with stratification stem from the non-zero stratum.

## 239 **2.7 Field Data Analysis**

240 HabCam images from the US sea scallop survey (NEFSC, 2014) were used to illustrate  
241 the usefulness of the methods discussed above on real data. For testing purposes, all the  
242 images were annotated using computer vision software (Gallager et al., unpublished) and  
243 also manually annotated, so that the estimates can be compared to their true values.

244 The automated annotator used a series of features including texture, color, and shape.  
245 A kernel of 100 x 100 pixels was run through each image left to right, top to bottom,  
246 extracting each feature set resulting in a feature vector of length 480 by width 3 (texture,  
247 color, and shape). Texture features were extracted using a 2-dimensional Gabor wavelet  
248 convolved with Gaussian kernels at 360 orientations for each pixel box providing rotational  
249 independent texture features (Gallager and Tiwari, 2008). Color was extracted in L\*A\*B\*  
250 color space using the color angle approach, where the standard deviation of the gradient  
251 between the pixel radius at 10 degree increments was extracted with 128 colors (Gallager  
252 and Tiwari, 2008). For each kernel, a Canny edge detection algorithm was used followed by  
253 extraction of Fourier shape descriptors. A Principal Component Analysis was run to reduce  
254 data dimensionality from  $> 4000$  to 128 principal components. Finally, a linear Support  
255 Vector machine was trained on 3800 images containing scallops of various sizes as well as

256 images containing no scallops over varying substrate conditions. The result was a probability  
 257 of the presence of a scallop; a scallop was considered as detected if this probability was greater  
 258 than 90%.

259 One out of every 50 images collected were annotated manually as well as with software  
 260 (Table 1), and this collection of images served as the data for our analysis. Data from  
 261 three regions with various probability of a false negative were selected. The probability of  
 262 a false positive could not be defined for our datasets because number of possible distractors  
 263 for each image was not identified. For each region, the manual annotations from a 7%  
 264 random subset of the images were used for estimation along with automated annotations  
 265 from each image; error rates could therefore be assessed because each image in the datasets  
 266 were annotated manually, even though only a sample of the manual annotations were used  
 267 in the analysis. The manual annotation subset was resampled 2000 times, and the various  
 268 estimation methods were applied to each iteration.

269 In the field, factors such as vehicle altitude, depth, etc. may also influence the performance  
 270 of the estimators. We tested an additional method that included auxiliary variables in the  
 271 two-strata local regression:

$$Y_j = a_0(u_j, v_j) + a_1(u_j, v_j)X_j + \sum_{b=2}^5 a_b(u_j, v_j)A_{bj} + \epsilon_j \quad (12)$$

272 where  $(u_j, v_j)$  is the coordinates of point  $j$  and  $a_b(u_j, v_j)$ 's are the coefficients of variables  $A$   
 273 including altitude, depth, squared depth, and latitude at location  $(u_j, v_j)$  point  $j$ .

## 274 2.8 Evaluation of Methods

275 For both simulation and field data analysis, mean squared error (MSE) and mean absolute  
 276 error (MAE) were used as the principal measures of precision and bias:

$$\text{MSE} = \frac{1}{K} \sqrt{\sum_{k=1}^K (Z_k - \mu)^2} \text{ and } \text{MAE} = \frac{1}{K} \sum_{k=1}^K |Z_k - \mu|, \quad (13)$$

277 where  $Z$  is the population estimates based on automated and manual annotations,  $\mu$  is the  
 278 true population abundance, and  $K$  is the number of iterations. These were reported relative

279 to the global unstratified manual sample mean (M1G):

$$\text{MSE}_{\text{re}} = \frac{\text{MSE} - \text{MSE}_{\text{M1G}}}{\text{MSE}_{\text{M1G}}} \text{ and } \text{MAE}_{\text{re}} = \frac{\text{MAE} - \text{MAE}_{\text{M1G}}}{\text{MAE}_{\text{M1G}}}. \quad (14)$$

280 MAE and MSE both reflect precision as well as bias but MSE weights more on large errors  
281 than small ones.

## 282 **3 Results**

### 283 **3.1 Simulation Results**

284 Combining automated and manual annotations using our methods increased precision of the  
285 estimates over manual counts alone by up to a maximum of 73%, whereas using the uncor-  
286 rected automated counts could decrease both accuracy and precision up to 717%, compared  
287 to using manual counts only (Tables 2 and 3). Increasing the number of manual samples  
288 increased the precision of all methods, but only by a modest amount (up to 15%).

289 In the base case, splitting the annotations into two strata was the most effective way of  
290 improving estimation precision, except at very high false negative rates where stratification  
291 degraded the estimates (Figures 3 and 4). When both false negative and false positive rates  
292 are low, the use of automated data for stratification and/or estimation substantially improves  
293 the precision and accuracy of the estimates regardless of the estimator used. Local models  
294 were superior to global models only when stratification was employed. For one-stratum  
295 allocation and when the probability of a false positive is high, the ratio and regression  
296 estimator performed better, whereas the offset estimator was better when the probability  
297 of a false negative is high but false positive rate is low. Similar patterns were observed for  
298 the other scenarios tested, i.e., the performance of the bias correction methods we tested are  
299 robust to changes in the environment and species distributions.

300 The probability of a false negative is the key factor determining the most effective bias  
301 correction methods, regardless of the level of probability of a false positive (Tables 2-5).  
302 When the probability of a false negative is low, nearly all the methods can improve the  
303 accuracy and precision of the population estimates, but stratification with the local ratio or  
304 the local regression estimator was generally superior. If the probability of false negatives is

305 high, no stratification with a simple global linear regression or manual sampling alone tended  
306 to have the best performance. If in addition the false positive rate is low, the global offset  
307 estimator also performs well.

### 308 **3.2 Field Data Analysis Results**

309 Results from the field data analysis were consistent with those from the simulations. Esti-  
310 mations of the mean using automated annotations alone were 63% to 498% higher than the  
311 simple manual sample mean (Table 1). For the region with low false negative rates (0.31), the  
312 two-strata local regression without auxiliary variables and two-strata local ratio estimator  
313 were superior; these increased precision over the simple manual sample mean by up to 51%  
314 (Table 1). When the false negative rate was higher (0.73-0.75), global regression or simply  
315 the manual sample mean were the best, with the global regression model improving precision  
316 by at most 11% over the simple manual sample mean. The offset estimator performed better  
317 than the ratio estimator in one case, likely because the false positive rate of this dataset is  
318 low; however, this is not totally clear since the false positive rates were not available for all  
319 of our field data. Auxiliary variables did not improve the performance of local regression for  
320 these data.

## 321 **4 Discussion**

322 The results indicate that combining even a mediocre automated annotator with manual  
323 annotations may be able to improve statistical efficiency over manual annotations alone when  
324 using the methods presented here. The combination of automated and manual annotations  
325 outperformed manual or (unadjusted) automated annotations alone, even when the false  
326 positive and false negative rates were as high as 0.5. The results from both simulations and  
327 field data analysis are consistent, and indicate that probability of a false negative is the  
328 key factor determining the best estimation method. The probability of a false positive does  
329 matter to some extent, especially when the probability of a false negative is higher, but even  
330 in this case, it is not the main factor determining the best method.

331 Stratification based on zero and positive automated counts is the most effective technique

332 to improve precision except at very high false negative rates. Stratification directly improves  
333 precision when the within-strata variance is less than the between strata variance (Cochran,  
334 1977), which is likely to be the case for even a moderately effective automated annotator.  
335 In addition, the allocation of manual samples between the two strata often further increases  
336 performance by allocating disproportionately more manual samples to the more variable  
337 stratum. Stratified estimates are in particular more precise at high false positive but low  
338 false negative rates. The zero stratum has no false positives, and contains a limited number  
339 of actual targets when the false negative rates are low. The zero stratum thus tends to have  
340 a low variance, so the number of targets in this stratum can be estimated precisely by a  
341 relatively small number of manual samples. This allows for higher sampling rates in the  
342 non-zero stratum, increasing the precision there.

343 The simple two-strata stratification presented here is natural for zero-inflated data such  
344 as in our examples. In some cases, more complex stratification may give further benefits. For  
345 example, there could be three strata, composed of where the automated annotator detects  
346 zero, one or more than one targets. We implicitly assumed for simplicity that the cost of a  
347 manual annotation is the same in each stratum. In reality, the labor cost of annotating an  
348 image tends to go up with the number of targets in the image. If this cost function is known,  
349 it can be taken into account in the optimal allocation among strata (Cochran, 1977).

350 In real world situations, the false negative (and positive) rates may be uncertain. In such  
351 cases, we recommend manually annotating a small sample of images to roughly estimate  
352 this rate, and select the manual sampling strategy (e.g., stratification scheme) and estimator  
353 based on this information. The optimal strategy is fairly robust to modest changes in the  
354 automated annotator error rates, so only a crude estimate of the false negative rates is needed  
355 to design a sampling strategy.

356 The offset estimator, by its definition, can account for errors that are independent of the  
357 target density, but less efficient in tracking errors that vary with the targets. Conversely, the  
358 ratio estimator is more effective without stratification when there are false negatives but few  
359 false positives (Figures 3 and 4), because the ratio estimator can take into account errors  
360 that are proportional to the target density. The precision of the ratio estimator depends on  
361 the correlation between automated and true counts (see Appendix); false positives directly

362 reduce this correlation.

363 In principle, the regression estimator should be able to account for both these types of  
364 errors, but it has the disadvantage of having two parameters that can be confounded with  
365 each other, especially at low sample sizes and when the data are zero-inflated. For stratified  
366 local regressions, the manual sample size used to estimate the regression parameters at each  
367 location is low, and might be one of the reasons why its performance is slightly lower than the  
368 stratified local ratio estimator. The difference in performance of stratified local regression  
369 estimator and stratified local ratio estimator was larger when the manual sample size is only  
370 1% and became smaller as the manual sample size increased (Tables 2 and 3).

371 There are nonetheless some advantages of regression methods. For example, multiple  
372 regression can be used if there is more than one automated annotator available, using counts  
373 from each automated annotator as predictors. Even though in our example field data it  
374 was not effective, auxiliary variables such as water depth, latitude, or substrate type may  
375 sometimes also be useful as predictors in a multiple regression.

376 Local estimation methods can improve estimates when the distribution of targets or  
377 errors is autocorrelated. In particular, false positives induced by distractors such as sand  
378 dollars and dead scallop shells are typically autocorrelated. False negative rates could be in  
379 some cases also autocorrelated (caused by e.g., poor visibility), but this would normally be  
380 a weaker effect than false positives if it exists at all. Stratification isolates the false positives  
381 in one stratum, which may be the reason that it enhances the effectiveness of using local  
382 estimation methods. The benefits of local estimation methods are however minor compared  
383 to stratification, even in the presence of substantial autocorrelation.

384 Although computer vision methods are rapidly improving, it is unlikely that automated  
385 detection of underwater organisms will be error free in the foreseeable future. Many marine  
386 organisms are cryptic, and can adjust their pattern and coloration to match their surround-  
387 ings, thus making it difficult to totally eliminate false negatives. For scallops in particular,  
388 false negatives can be caused by colonization of their shell by epifauna or the shell being  
389 covered by marine snow or sediments. In addition, a small percentage ( $\sim 5-10\%$ ) of sea  
390 scallops are “albinos”, with white upper shells, that are difficult to distinguish from dead  
391 scallop shells. While we believe that the false positives induced by sand dollars can be



392 reduced considerably compared to present methods, it is also unlikely that false positives  
393 can be completely eliminated (for example, it is sometimes difficult to distinguish a dead  
394 scallop shell from a live scallop). Thus, combining automated and manual annotations using  
395 the methods described here is likely to continue to be an improvement over using either  
396 automated or manual annotations alone.

397 While we have focused on automated annotations of marine organisms, our methods are  
398 applicable to a much wider set of problems. For example, our methods could be employed  
399 whenever there are at least two observers counting the same things, one of whom is an  
400 expert (or is a reference collection) who is considered error free but only observes a sample.  
401 Annotations using crowd-sourcing (Simpson et al., 2014) may be subject to higher error rates  
402 than those done by experts, which can be corrected using the techniques presented here.  
403 Our methods also are applicable to automated or crowd-sourced annotations of a variety  
404 of targets beyond those underwater, such as targets from aerial photography, surveillance  
405 cameras, medical imaging and testing, and industrial quality control.

## 406 **5 Acknowledgements**

407 We thank Chris Legault and two anonymous reviewers for their constructive comments.  
408 The research was supported in part by funding from the NOAA Automated Image Analysis  
409 Strategic Initiative.

## 410 **References**

411 Beijbom, O. 2014. Random Sampling in an Age of Automation: Minimizing Expenditures  
412 through Balanced Collection and Annotation. arXiv preprint arXiv:1410.7074.

413

414 Beijbom, O., Edmunds, P.J., Kline, D., Mitchell, B.G., Kriegman, D., 2012. Automated  
415 annotation of coral reef survey images. IEEE Conference In Computer Vision and Pattern  
416 Recognition (CVPR), p. 1170-1177.

417

418 Brand, A.R., 2006. Scallop ecology: distributions and behavior. In: *Scallops: Biology, Ecology, and Aquaculture*, Elsevier, Amsterdam, p. 651-744.

419

420

421 Brunson, C., Fotheringham, A.S., and Charlton, M., 2008. Geographically weighted regression: a method for exploring spatial nonstationarity. *Encyclopedia of Geographic Information Science*, 558.

422

423

424

425 Culverhouse, P.F., Williams, R., Benfield, M., Flood, P. R., Sell, A.F., Mazzocchi, M.G., Buttino, I. Sieracki, M., 2006. Automatic image analysis of plankton: Future perspectives. *Mar. Ecol. Prog. Ser.* 312, 297-309.

426

427

428

429 Cochran, W.G., 1977. *Sampling Techniques: 3rd Ed.*, Wiley, New York.

430

431 Cressie, N.A.C., 1993. *Statistics for Spatial Data*, revised edition. Wiley, New York.

432

433 Davis, C.S., Gallagher, S.M., Solow, A.R., 1992. Microaggregations of oceanic plankton observed by towed video microscopy. *Science* 257(5067), 230-232.

434

435

436 Dawkins, M., Stewart, C., Gallagher, S., York, A., 2013. Automatic scallop detection in benthic environments, 2013 IEEE Workshop on Applications of Computer Vision (WACV), p. 160-167.

437

438

439

440 Gallagher, S.M., Singh, H., Tiwari, S., Howland, J., Rago, P., Overholtz, W., Taylor, R., Vine, N., 2005. High resolution underwater imaging and image processing for identifying essential fish habitat. Report of the National Marine Fisheries Service Workshop on Underwater Video analysis. DA Somerton and CT Glendill (eds) NOAA Technical Memorandum NMFS-F/SPO-68. p. 44-54.

441

442

443

444

445

446 Gallagher, S.M., Tiwari, S., 2008. Optical method and system for rapid identification of multiple refractive index materials using multiscale texture and color invariants. US patent

447

448 Number 7,415,136.

449

450 Gallagher, S.M., Nordahl, V., Godlewski, J.M., 2014. The habitat mapping camera system  
451 (HabCam). Proceedings of the Undersea Imaging Workshop. R Langton and P Rowe, Eds.

452

453 Gallagher, S.M., Honig, P., York, A.D, Hart D.R., Unpublished. Automated detection and  
454 classification of benthic fauna along the US Northeast Continental Shelf. Unpublished Re-  
455 sults.

456

457 Guo, L., Ma, Z., Zhang, L., 2008. Comparison of bandwidth selection in application of geo-  
458 graphically weighted regression: a case study. Can. J. Fish. Aquat. Sci. 38(9), 2526-2534.

459

460 Hart, D.R., 2006. Effects of sea stars and crabs on sea scallop (*Placopecten magellanicus*)  
461 recruitment in the Mid-Atlantic Bight. Mar. Ecol. Prog. Ser. 306, 209-221.

462

463 Howland, J., Gallagher, S.M., Singh Girard, H., Abrams, L., Griner, C., 2006. Development  
464 of a towed, ocean bottom survey camera system for deployment by the fishing industry.  
465 IEEE Oceans p. 1-10.

466

467 Hu, Q., Davis, C.S., 2006. Accurate automatic quantification of taxa-specific plankton abun-  
468 dance using dual classification with correction. Mar. Ecol. Prog. Ser. 306, 51-61.

469

470 Kannappan, P., Walker, J.H., Trembanis, A., Tanner, H.G., 2014. Identifying sea scallops  
471 from benthic camera images. Limnol. Oceanogr.: Methods 12, 680–693.

472

473 Marcos, M.S.A., David, L., Peñaflor, E., Ticzon, V., Soriano, M., 2008. Automated benthic  
474 counting of living and non-living components in Ngedarrak Reef, Palau via subsurface un-  
475 derwater video. Environ. Monit. Assess. 145(1-3), 177-184.

476

477 Northeast Fisheries Science Center [NEFSC]. 2014. 59th Northeast Regional Stock Assess-

478 ment Workshop: Assessment Report. Northeast Fisheries Science Center Reference Docu-  
479 ment 14-09.

480

481 Rosenkranz, G.E., Gallager, S.M., Shepard, R.W., Blakesleed, M., 2008. Development of a  
482 high-speed, megapixel benthic imaging system for coastal fisheries research in Alaska. *Fish.*  
483 *Res.* 92, 340–344.

484

485 Simpson, R., Page, K.R., De Roure, D., 2014. Zooniverse: observing the world’s largest  
486 citizen science platform. In: Proceedings of the companion publication of the 23rd inter-  
487 national conference on the World Wide Web International World Wide Web Conferences  
488 Steering Committee, p. 1049-1054.

489

490 Singh, W., Örnólfssdóttir, E.B., Stefansson, G., 2013. A camera-based autonomous underwa-  
491 ter vehicle sampling approach to quantify scallop abundance. *J. Shellfish Res.* 32(3),725-732.

492

493 Solow, A., Davis C., Hu, Q., 2001. Estimating the taxonomic composition of a sample when  
494 individuals are classified with error. *Mar. Ecol. Prog. Ser.* 216, 309-311.

495

496 Spampinato, C., Chen-Burger, Y.H., Nadarajan, G., Fisher, R.B. 2008. Detecting, tracking  
497 and counting fish in low quality unconstrained underwater videos. *VISAPP (2)*, 2008, 514-  
498 519.

499

500 Taylor, R., Vine, N., York, A., Lerner, S., Hart, D., Howland, J., Prasad, L., Mayer, L.,  
501 Gallager, S., 2008. Evolution of a benthic imaging system from a towed camera to an auto-  
502 mated habitat characterization system. *IEEE Oceans* p. 1-7.

503

504 Tolimieri, N., Clarke, M.E., Singh, H., Goldfinger, C. 2008. In: Reynolds, J.R. and Greene,  
505 H.G. (eds.), 2008. *Marine Habitat Mapping Technology for Alaska*, Alaska Sea Grant Col-  
506 lege Program, University of Alaska Fairbanks.

507

508 Verikas, A., Gelzinis, A., Bacauskiene, M., Olenina, I., Vaiciukynas, E., 2015. An Integrated  
509 Approach to Analysis of Phytoplankton Images. IEEE J. Ocean. Eng. 40(2), 315-326.

510

511 Yoklavich, M.M., Love, M.S., Forney, K.A. 2007. A fishery-independent assessment of an  
512 overfished rockfish stock, cowcod (*Sebastes levis*), using direct observations from an occupied  
513 submersible. Can. J. Fish. Aquat. Sci. 64(12), 1795-1804.

514

Table 1: Relative mean squared error ( $MSE_{Re}$ ) and relative mean absolute error ( $MAE_{Re}$ ) for each estimator, using unstratified (one-stratum) or two strata estimation, and either local or global estimation for three sets of actual HabCam field data. Error rates are relative to the global unstratified manual mean, which is used as a baseline. “AUTO” represents  $MSE_{Re}$  or  $MAE_{Re}$  calculated using only the automated annotations. “L+Var” represents local regression with auxiliary variables. The dark and light grey-shaded entries represent the best and second best method, respectively.

Sample Size	False Negative	Stat	Auto	Manual Mean				Ratio Est.				Offset Est.				Regression Est.				
				One-stratum		Two-strata		One-stratum		Two-strata		One-stratum		Two-strata		One-stratum		Two-strata		L+Var
				Global	Local	Global	Local	Global	Local	Global	Local	Global	Local	Global	Local	Global	Local	Global	Local	L+Var
5057	0.31	$MSE_{Re}$	1.68	0	-0.06	-0.04	-0.21	-0.07	-0.31	-0.09	-0.50	-0.02	0.01	-0.02	-0.19	-0.07	-0.26	-0.10	-0.51	-0.04
		$MAE_{Re}$	0.78	0	-0.03	-0.02	-0.11	-0.04	-0.16	-0.04	-0.29	-0.01	0.01	-0.00	-0.10	-0.04	-0.14	-0.05	-0.31	-0.08
9610	0.73	$MSE_{Re}$	4.98	0	0.04	-0.04	0.02	-0.06	-0.01	-0.03	0.03	-0.10	-0.05	-0.04	0.02	-0.11	-0.07	-0.03	0.03	0.05
		$MAE_{Re}$	1.68	0	0.01	-0.02	0.01	-0.03	-0.01	-0.01	0.01	-0.06	-0.03	-0.02	0.01	-0.06	-0.04	-0.02	0.01	0.01
14856	0.75	$MSE_{Re}$	1.25	0	0.71	0.37	0.16	0.28	2.06	0.89	0.55	0.40	1.90	1.16	1.61	-0.00	0.93	0.37	0.07	0.04
		$MAE_{Re}$	0.63	0	0.36	0.17	0.07	0.13	0.88	0.37	0.26	0.18	0.75	0.47	0.62	-0.00	0.47	0.17	0.03	0.02







Table 4: Proportion of runs with the least mean square error (MSE) for the five scenarios by type of estimators. See Table 2 for explanations of the notations.

Scenarios	S2: Good Vis.	0.16	0	0	0	0	0.09	0.08	0.03	0.32	0.01	0	0.01	0	0.05	0	0.07	0.18	S1: F1S<=0.5 & F1D<=0.5
	Mod.	0.15	0	0	0	0	0.11	0.07	0.05	0.29	0.01	0	0	0.01	0.06	0	0.08	0.17	
	Poor	0.13	0	0	0	0	0.1	0.07	0.03	0.26	0.02	0	0.02	0.02	0.05	0	0.12	0.18	
	S3: Neg. Cor.	0.09	0	0	0.01	0.01	0	0.02	0	0.51	0	0	0	0	0	0.03	0	0.32	
	Pos.	0.07	0.02	0	0.02	0	0	0	0	0.74	0	0	0	0	0	0.01	0.02	0.13	
	None	0.07	0.01	0	0.02	0	0	0.01	0	0.45	0	0	0	0	0.02	0.01	0	0.41	
	S4: High Autocor.	0.1	0	0	0	0	0.08	0.07	0.02	0.37	0	0	0.01	0	0.02	0	0.05	0.28	
	Medium	0.16	0	0	0	0	0.09	0.08	0.03	0.32	0.01	0	0.01	0	0.05	0	0.07	0.18	
	Low	0.15	0	0	0.02	0	0.16	0.01	0.07	0.28	0.02	0.01	0.02	0.04	0.07	0	0.12	0.03	
	S5: M/T 1%	0.4	0	0	0	0	0.07	0	0.06	0.16	0.02	0	0.07	0.01	0.06	0.01	0.09	0.06	
	3%	0.22	0	0	0	0	0.14	0.04	0.04	0.23	0.04	0	0.02	0.02	0.04	0	0.05	0.15	
	7%	0.16	0	0	0	0	0.09	0.08	0.03	0.32	0.01	0	0.01	0	0.05	0	0.07	0.18	
	11%	0.11	0	0	0	0	0.09	0.04	0.02	0.32	0.02	0	0.01	0.01	0.04	0	0.12	0.22	
	15%	0.1	0	0	0.02	0	0.11	0.06	0.03	0.31	0	0	0	0.02	0.03	0	0.09	0.23	
	S2: Good Vis.	0	0	0	0.01	0.01	0	0.06	0	0.61	0	0	0	0	0.02	0.01	0	0.26	
Mod.	0	0.01	0	0.02	0	0	0.05	0	0.64	0	0	0	0	0.01	0.02	0	0.26		
Poor	0	0.03	0	0.02	0	0	0.02	0	0.68	0	0	0	0	0.03	0.03	0	0.2		
S3: Neg. Cor.	0	0	0	0.01	0.02	0	0.01	0	0.56	0	0	0	0	0.01	0	0	0.39		
Pos.	0	0	0	0.08	0.04	0	0	0	0.58	0	0	0	0	0.01	0	0.09	0.2		
None	0	0	0	0.01	0.01	0	0.01	0	0.57	0	0	0	0	0.01	0.01	0	0.38		
S4: High Autocor.	0	0	0	0	0	0	0	0	0.56	0	0	0	0	0	0.01	0	0.43		
Medium	0	0	0	0.01	0.01	0	0.06	0	0.61	0	0	0	0	0.02	0.01	0	0.26		
Low	0	0.02	0	0.08	0.03	0.05	0.01	0	0.67	0	0	0	0	0.07	0.01	0.04	0.03		
S5: M/T 1%	0.02	0.05	0	0.25	0.02	0.04	0.02	0.09	0.34	0	0	0	0	0.01	0	0.02	0.16		
3%	0	0.03	0	0.08	0.02	0.05	0.04	0	0.5	0	0	0	0	0.02	0.05	0.01	0.21		
7%	0	0	0	0.01	0.01	0	0.06	0	0.61	0	0	0	0	0.02	0.01	0	0.26		
11%	0	0	0	0	0	0	0	0	0.68	0	0	0	0	0	0.02	0	0.3		
15%	0	0.01	0	0	0	0	0.02	0	0.68	0	0	0	0	0.01	0.01	0	0.27		
S2: Good Vis.	0.2	0.1	0	0	0	0.09	0.07	0	0	0.19	0	0	0	0.34	0.01	0	0	S1: F1S>0.5 & F1D<=0.5	
Mod.	0.25	0.03	0	0	0	0.09	0.08	0	0	0.24	0	0	0	0.3	0.01	0	0		
Poor	0.25	0.05	0	0	0	0.09	0.08	0	0	0.24	0	0	0	0.28	0	0	0		
S3: Neg. Cor.	0.47	0.25	0	0.05	0.02	0	0	0.02	0	0.02	0	0	0	0.08	0.05	0.04	0.01		
Pos.	0.31	0.23	0	0.17	0.05	0	0.01	0.03	0	0.01	0	0	0	0.09	0.03	0.05	0.03		
None	0.44	0.24	0	0.03	0.02	0	0.01	0.03	0	0.01	0	0	0	0.09	0.05	0.05	0.04		
S4: High Autocor.	0.16	0.05	0	0.04	0	0.06	0.25	0	0	0.05	0	0	0	0.18	0.18	0.01	0		
Medium	0.2	0.1	0	0	0	0.09	0.07	0	0	0.19	0	0	0	0.34	0.01	0	0		
Low	0.15	0.12	0	0.03	0	0.11	0.01	0.01	0	0.21	0	0.01	0	0.35	0	0.02	0		
S5: M/T 1%	0.45	0.08	0	0.04	0	0.06	0.01	0.01	0	0.14	0	0	0.03	0.15	0.04	0	0		
3%	0.25	0.15	0	0	0.01	0.12	0.06	0.02	0	0.15	0	0.01	0	0.21	0.01	0.01	0		
7%	0.2	0.1	0	0	0	0.09	0.07	0	0	0.19	0	0	0	0.34	0.01	0	0		
11%	0.11	0.16	0	0.01	0	0.09	0.13	0.02	0	0.16	0	0	0	0.25	0.05	0	0.01		
15%	0.11	0.12	0	0.05	0	0.15	0.19	0.01	0	0.08	0	0.01	0	0.18	0.1	0	0.01		
S2: Good Vis.	0.21	0.25	0	0.04	0.03	0	0.1	0.02	0.02	0	0	0	0	0.24	0.04	0.03	0.02		S1: F1S>0.5 & F1D>0.5
Mod.	0.2	0.31	0	0.01	0	0.02	0.1	0	0.01	0	0	0	0	0.21	0.09	0.03	0.02		
Poor	0.19	0.23	0	0.04	0.02	0	0.11	0	0.01	0	0	0	0	0.3	0.03	0.03	0.04		
S3: Neg. Cor.	0.09	0.31	0	0.08	0.05	0	0.08	0.01	0.03	0	0	0.01	0	0.18	0.01	0.1	0.05		
Pos.	0	0.24	0	0.14	0.2	0	0	0.03	0.01	0	0	0	0	0.16	0	0.2	0.02		
None	0.04	0.32	0	0.11	0.05	0	0.08	0.01	0.03	0	0	0	0	0.16	0.04	0.12	0.04		
S4: High Autocor.	0.14	0.08	0	0.01	0.01	0	0.05	0.02	0.05	0	0	0	0	0.05	0.44	0.01	0.14		
Medium	0.21	0.25	0	0.04	0.03	0	0.1	0.02	0.02	0	0	0	0	0.24	0.04	0.03	0.02		
Low	0.31	0.26	0	0.11	0.02	0.04	0	0.01	0	0	0	0	0	0.17	0	0.07	0.01		
S5: M/T 1%	0.64	0.18	0.01	0.08	0.02	0.02	0	0	0.01	0	0	0	0	0.03	0	0.01	0		
3%	0.37	0.29	0	0.07	0.02	0	0.01	0	0	0	0	0	0	0.12	0.03	0.07	0.02		
7%	0.21	0.25	0	0.04	0.03	0	0.1	0.02	0.02	0	0	0	0	0.24	0.04	0.03	0.02		
11%	0.17	0.25	0	0.03	0.02	0	0.1	0.03	0.07	0	0	0	0	0.22	0.04	0	0.07		
15%	0.15	0.16	0	0.03	0.04	0	0.18	0.04	0.04	0	0	0	0	0.17	0.09	0.04	0.06		
		AUTO	M1G	M1L	M2G	M2L	Ra1G	Ra1L	Ra2G	Ra2L	O1G	O1L	O2G	O2L	Re1G	Re1L	Re2G	Re2L	



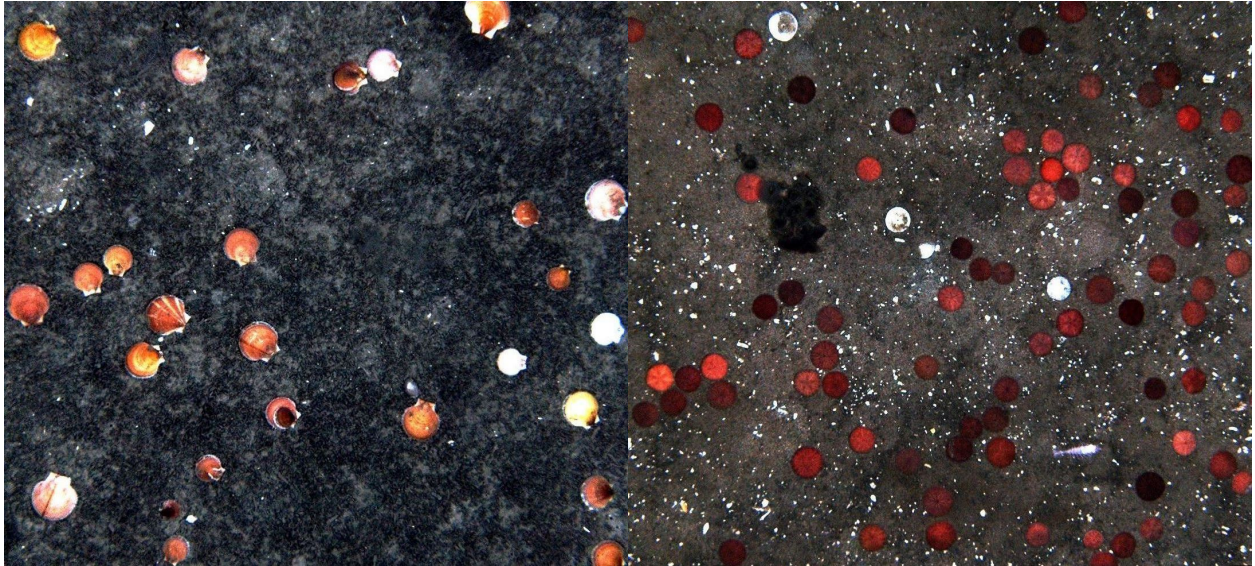


Figure 1: HabCam Images with scallops (left) and its common distractor sand dollars (right).

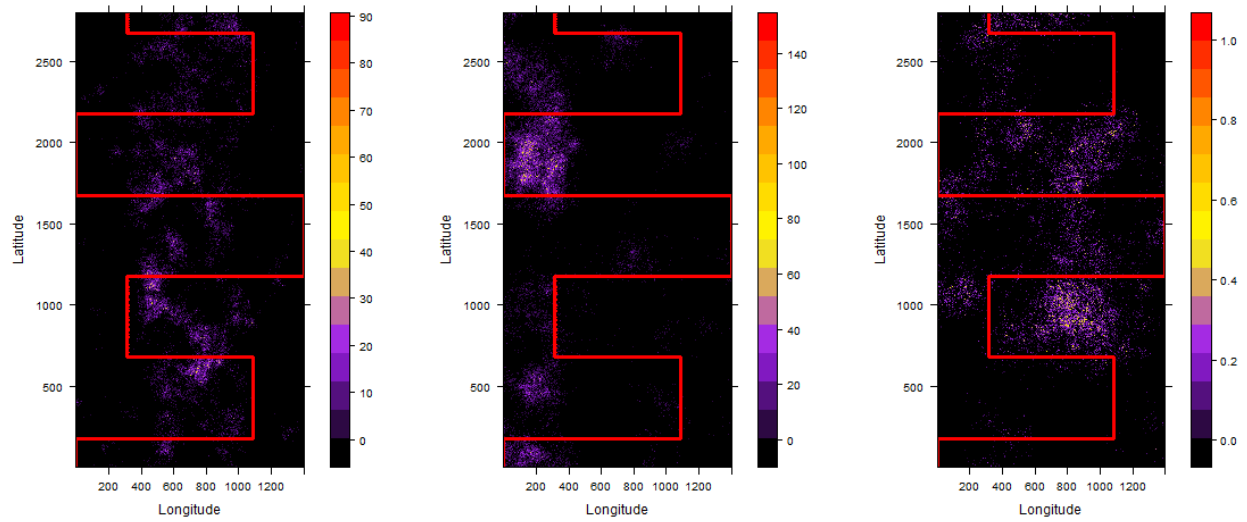


Figure 2: Example simulated distributions of scallops (left), distractors (center; moderate autocorrelation and negatively correlated with scallop distribution), and water visibility (right; poor) with an over-layed sampling track (red line). The colors represent counts per m<sup>2</sup> for scallops and distractors and the reduced probabilities of detecting scallops and distractors due to poor water visibility.

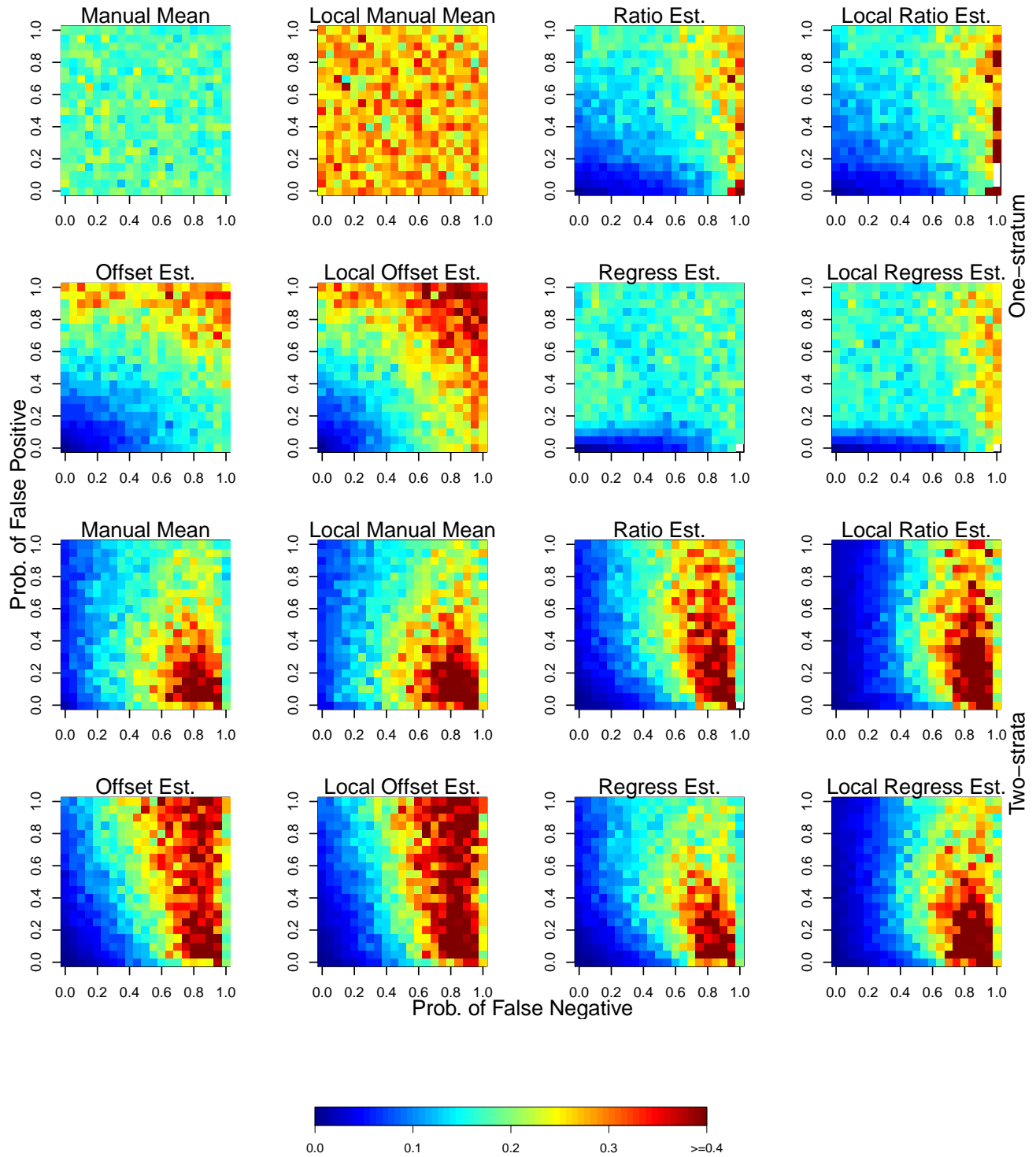


Figure 3: Mean squared error (MSE, indicated by color) at various false negative and false positive rates in the base case scenario, by estimator type, global or local estimation, and unstratified (one-statum) or two-strata estimation.

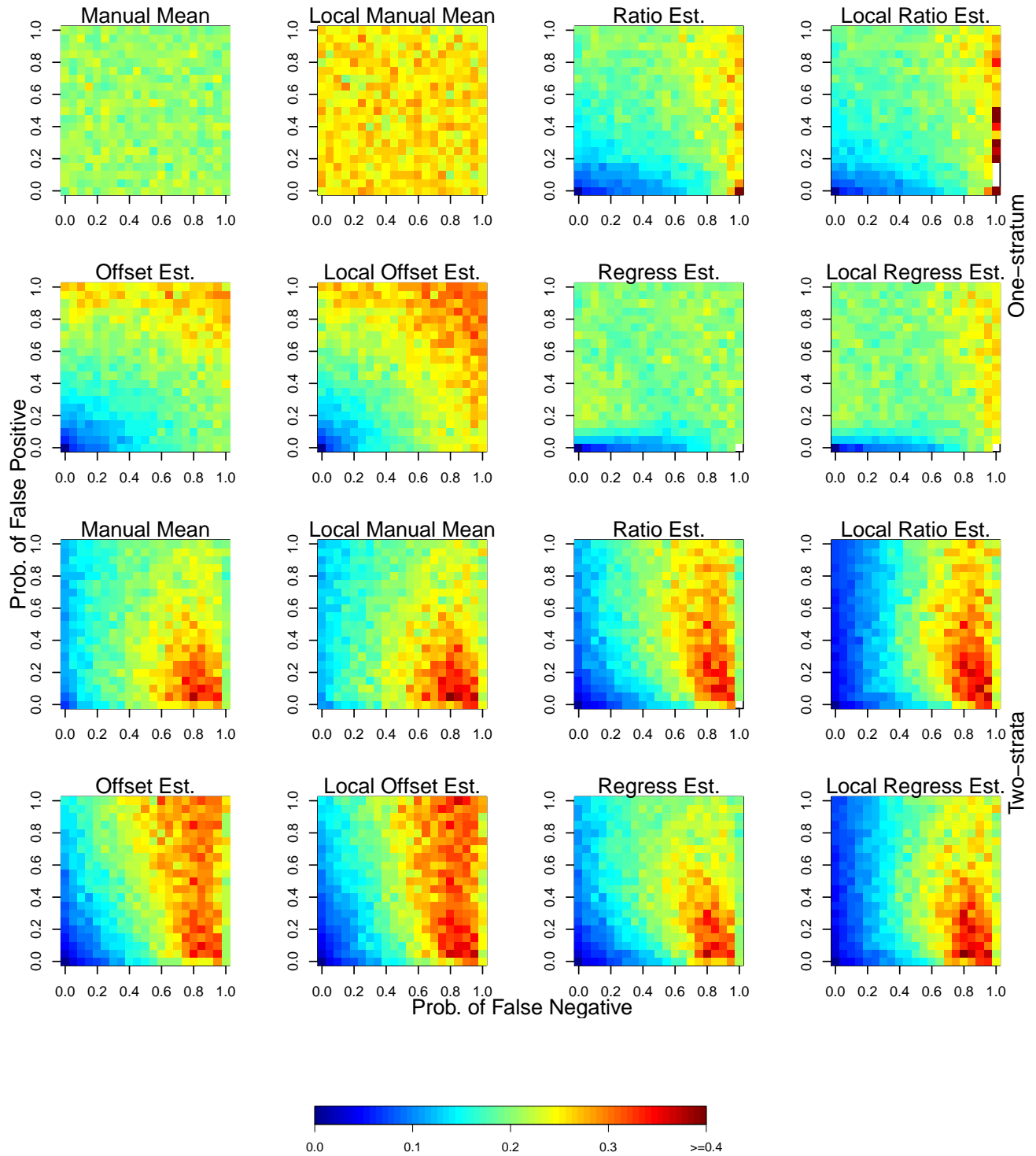


Figure 4: Mean absolute error (MAE, indicated by color) at various false negative and false positive rates in the base case scenario, by estimator type, global or local estimation, and unstratified (one-stratum) or two-strata estimation.

515 **Appendix - Analytic derivation of properties of the ratio estimator**

516 Let  $Y_i$  be the number of targets in the  $i$ th randomly chosen image; it will be assumed that  
 517 manual processing is perfect, so that  $Y_i$  is also the number of targets that were detected  
 518 manually. Let  $X_i$  be the number of targets detected by the automated software in the  $i$ th  
 519 image. We will consider the following ratio estimator for the mean number of targets:

$$T = \mu_X \frac{Y_1 + Y_2 + \dots + Y_n}{X_1 + X_2 + \dots + X_n} = \mu_X \frac{\bar{Y}}{\bar{X}} \quad (15)$$

520 where  $\mu_X$  is the mean of the automated counts over all photographs, and  $\bar{X}$  and  $\bar{Y}$  are the  
 521 sample means for the automated and manual counts for a randomly chosen sample of  $n$   
 522 images. Let  $\mu_X = E(X_i)$  and  $\mu_Y = E(Y_i)$ ,  $\sigma_X$  and  $\sigma_Y$  be the standard deviations of  $X_i$  and  
 523  $Y_i$ , respectively, and let  $\rho$  be the correlation between  $X_i$  and  $Y_i$ . Assuming for simplicity  
 524 that the finite population correction factor is negligible (i.e., that the total number of images  
 525 is large relative to  $n$ ; this does not affect the main results below), using the approximate  
 526 variance for a ratio (Cochran, 1977),

$$\text{Var}(T) = \mu_X^2 \text{Var} \frac{\bar{Y}}{\bar{X}} \simeq \mu_X^2 \frac{1}{\mu_X^2} \left[ \sigma_Y^2 + \frac{\sigma_X^2 \mu_Y^2}{\mu_X^2} - 2\rho \sigma_X \sigma_Y \frac{\mu_Y}{\mu_X} \right] / n \quad (16)$$

$$= \left[ \sigma_Y^2 + \sigma_X \frac{\mu_Y}{\mu_X} (\sigma_X \frac{\mu_Y}{\mu_X} - 2\rho \sigma_Y) \right] / n. \quad (17)$$

527 Hence,  $\text{Var}(T)$  decreases linearly with  $\rho$ . If  $\mu_X = \mu_Y$  and  $\sigma_X = \sigma_Y$ , this reduces to  $\text{Var}(T) \simeq$   
 528  $2\sigma_Y^2(1 - \rho)/n$ .

529 By comparison, a simple random sample of  $n$  manual images has variance  $\text{Var}(\bar{Y}) = \sigma_Y^2/n$ ,  
 530 which is the first term of equation (17). Thus, the ratio estimator  $T$  has lower variance than  
 531 simply using the manual images (i.e.,  $\text{Var}(T) < \text{Var}(\bar{Y})$ ) if and only if  $\sigma_X \frac{\mu_Y}{\mu_X} - 2\rho \sigma_Y < 0$ ,  
 532 i.e.,

$$\rho > \frac{\sigma_X \mu_Y}{2\sigma_Y \mu_X}. \quad (18)$$

533 In particular, if the  $X_i$ s and  $Y_i$ s have the same means and variances, then the ratio estimator  
 534 is an improvement over simple random sampling of the manual images if and only if  $\rho > 1/2$ .