

1 **Recruitment location influences bias and uncertainty in SARS-CoV-2 seroprevalence estimates**

2 Tyler S. Brown^{1,2,*}, Pablo Martinez de Salazar Munoz², Abhishek Bhatia², Bridget Bunda¹, Ellen K.
3 Williams, David Bor³, James S. Miller¹, Amir M. Mohareb¹, Julia Thierauf¹, Wenxin Yang¹, Julian
4 Villalba¹, Vivek Naranbai¹, Wilfredo Garcia Beltran¹, Tyler E. Miller¹, Doug Kress⁴, Kristen Stelljes⁵,
5 Keith Johnson⁵, Daniel B. Larremore⁶, Jochen Lennerz¹, A. John Iafrate¹, Satchit Balsari², Caroline O.
6 Buckee^{2,**}, Yonatan H. Grad^{2,**}

8 ¹ Massachusetts General Hospital, Boston, MA

9 ² Harvard T.H. Chan School of Public Health, Boston, MA

10 ³ Cambridge Health Alliance, Cambridge, MA

11 ⁴ Board of Health, City of Somerville, Massachusetts

12 ⁵ SomerStat, City of Somerville, Massachusetts

13 ⁶ University of Colorado, Boulder

14 * Corresponding author: tsbrown@mgh.harvard.edu

15 ** Joint senior authors

16 Abstract: 147 words

17 Main text: 3,225

18 **Abstract**

19 The initial phase of the COVID-19 pandemic in the US was marked by limited diagnostic testing,
20 resulting in the need for seroprevalence studies to estimate cumulative incidence and define epidemic
21 dynamics. In lieu of systematic representational surveillance, venue-based sampling was often used to
22 rapidly estimate a community's seroprevalence. However, biases and uncertainty due to site selection
23 and use of convenience samples are poorly understood. Using data from a SARS-CoV-2
24 serosurveillance study we performed in Somerville, Massachusetts, we found that the uncertainty in
25 seroprevalence estimates depends on how well sampling intensity matches the known or expected
26 geographic distribution of seropositive individuals in the study area. We use GPS-estimated foot traffic
27 to measure and account for these sources of bias. Our results demonstrated that study-site selection
28 informed by mobility patterns can markedly improve seroprevalence estimates. Such data should be
29 used in the design and interpretation of venue-based serosurveillance studies.

30 **Keywords:** SARS-CoV-2, serosurveillance, study design, human mobility

31 **Introduction**

32 Studies estimating SARS-CoV-2 seroprevalence have been critical to our understanding of COVID-19
33 pandemic dynamics, particularly in times when diagnostic testing was limited and the extent of
34 community spread was unknown [1–4]. However, many of these estimates were derived from non-
35 representational convenience sampling [4–6], thus making these estimates subject to multiple sources of
36 bias and uncertainty that remain poorly understood [7]. Low precision can obscure differences in
37 estimated seroprevalence between populations, limiting efforts to understand heterogeneity in epidemic
38 intensity or vaccination coverage and reducing the effectiveness of public health interventions. As cost,
39 speed, and logistical concerns continue to motivate the use of convenience sampling [8], identifying and
40 accounting for bias and uncertainty are therefore critical to improving the design and interpretability of
41 these studies.

42 Convenience sampling involves inherently non-uniform sampling across demographic or geographic
43 subgroups. For example, in venue-based (or “walk-up”) studies, in which participants are recruited from
44 among visitors to a central or highly trafficked location [4,9], the geographic distribution of participants
45 is expected to be skewed toward individuals living closer to the study location. A similar concern
46 applies to studies using discarded blood samples, where the catchment area of a given hospital or
47 clinical laboratory may strongly constrain the geographic distribution of samples available for analysis
48 [5].

49 Multiple recent studies have underscored the utility of GPS- and mobile phone-associated human
50 mobility data in understanding COVID-19 epidemiology. These studies have clarified the role of highly-
51 visited “super-spreader” locations [10,11] in local COVID-19 epidemics and identified higher aggregate
52 mobility as a predictor of neighborhood-level COVID-19 risk [3]. Prospective applications of these data
53 sources, including their use in the design and implementation of epidemiological studies, are limited.

54 Here, we examine seroprevalence estimates obtained via venue-based sampling [4,9], with direct
55 applications to other forms of convenience sampling [5,12,13]. We performed a venue-based
56 serosurveillance study in Somerville, Massachusetts and used these data to analyze bias and uncertainty
57 arising from inherent geographic variation in sampling intensity. This analysis showed that substantial
58 loss of precision can occur if the distribution of sampling intensity poorly matches the geographic
59 distribution of true seropositive individuals in the study area. As GPS-estimated foot traffic offers a
60 proxy measurement for the geographic distribution of visitors to different locations, we evaluated the
61 extent to which informed selection of locations for venue sampling reduces uncertainty and bias
62 introduced by geographic heterogeneity in underlying seropositivity. Our results thus offer an approach
63 to significantly improve the design and interpretation of seroprevalence studies that use venue-based
64 convenience sampling with little if any impact on cost and speed.

65 **Methods**

66 *SARS-CoV-2 seroprevalence study design and participant information*

67 We obtained serological and participant demographic data for 398 asymptomatic adults tested at a
68 temporary study site near an essential business location in Somerville, Massachusetts. The study was
69 conducted over 4 days (June 4th, 5th, 8th, and 9th, 2020), approximately 6 weeks after the first wave of the
70 COVID-19 epidemic peaked in Massachusetts [14]. The study was designated minimal risk human
71 subjects research and approved by institutional review boards at Massachusetts General Hospital and the
72 Harvard T.H. Chan School of Public Health (Protocol number: 2020P001081). The study recorded
73 participant demographic information including age, gender, and self-reported home locations (by ZIP
74 code and electoral ward). We also collected information on how participants learned about the study in
75 order to distinguish participants directly recruited on site at the study location from those who learned
76 about the study from friends, family, or social media. We did not advertise or announce enrollment for
77 the study prior to its implementation, with the goal of increasing the proportion of individuals recruited

78 on site at the study location. We used this data to calculate P_j^{direct} , the proportion of all directly recruited
79 participants from Somerville with self-reported home locations in electoral ward $j \in \{1, \dots, 7\}$. We refer
80 to $\{P_1^{\text{direct}}, \dots, P_7^{\text{direct}}\}$ as the “survey participant catchment distribution.”

81 *Public health acute infection data*

82 We obtained data on 916 PCR-confirmed COVID-19 cases with documented home addresses in
83 Somerville (collected from the onset of the epidemic through June, 2020) from the Massachusetts
84 Virtual Epidemiologic Network (MAVEN). COVID-19 cases are reported to MAVEN by state and local
85 health agencies, and cases are designated as “confirmed” if they have a positive result for SARS-CoV-2
86 RNA detection using an FDA-approved molecular amplification detection test, for example RT-PCR.
87 The total number of Somerville residents tested for SARS-COV-2 via RT-PCR was also obtained from
88 MAVEN. Data were anonymized and aggregated by electoral ward prior to analysis. We calculated the
89 cumulative incidence of PCR-confirmed infections by ward (θ_j^{PCR}) and the proportion of all PCR tests
90 with positive results (“PCR positivity”).

91 *GPS-estimated business foot traffic*

92 We used GPS-estimated foot traffic (SafeGraph, *safegraph.com*) to approximate the distribution of
93 home locations for daytime visitors for different locations of interest, which we refer to as the “GPS-
94 estimated catchment area.” We use these data to estimate home and work locations at the level of census
95 block group (CBG) for visitors to designated points-of-interest, such as businesses, and specific CBGs.
96 We used CBG-level visitor data for June 2020 as the primary data source for visitors in our analysis.
97 These data were filtered for CBGs with low visitor counts and re-aggregated from CBGs to electoral
98 wards prior to analysis (Supplementary Material, Figures S1 and S2). We used the filtered, re-
99 aggregated data to obtain GPS-estimated visitor catchment distributions for two locations: (1) the actual
100 study venue, located in Somerville electoral Ward 2 (denoted V_j^{site}) and (2) a hypothetical alternative
101 study venue located in Somerville electoral Ward 1 (denoted V_j^{alt}).

102 *Simulations*

103 We used numerical simulation to examine bias and uncertainty in estimated seroprevalence ($\hat{\theta}_{\text{pop}}$) as a
104 function of subgroup sizes, true subgroup seropositivity, and sample allocation. Here, “subgroup” is
105 used to describe any potential stratification of the overall population, including by demographic and/or
106 geographic characteristics. Briefly, we used demographic [15] and SARS-CoV-2 acute infection data
107 from Somerville, MA to generate a simulated population with varying true seropositivity $\theta_{j,k}$ across
108 subgroups stratified by age group j and location k (where locations are electoral wards in Somerville).
109 Briefly, the simulation randomly draws $n_{j,k}$ individuals from each subgroup, calculates weighted
110 population-level seroprevalence (adjusted for serological test performance), and repeats this process
111 10000 times to generate distributions of $\hat{\theta}_{\text{pop}}$ values. We report W , the width of the 95th percentile
112 interval for each distribution, as an approximate measure of uncertainty for $\hat{\theta}_{\text{pop}}$ for a given set of
113 simulation parameters. Additional details are available in the Supplementary Material. *R* code for the
114 numerical simulations is available at <https://github.com/svsero/COVID19serosurveillance-Somerville>.

115 **Results**

116 *Study participant catchment distributions and geographic heterogeneity in COVID-19 epidemic intensity*

117 We first examined how participant catchment distributions align with, or mismatch, the geographic
118 distribution of seropositive individuals in a given study area. We observed that the survey participant
119 catchment distribution in our serosurveillance study was skewed strongly toward locations near the
120 study site (Figure 1A). Among directly recruited participants with home locations in Somerville, 43%
121 (43/100) reported home locations in Somerville Ward 2 (where the study site was located) compared to
122 4% in Ward 1 and 4% in Ward 4. In contrast, the cumulative incidence of PCR-confirmed SARS-CoV-2
123 infections (θ_j^{PCR}) was approximately three-fold higher in electoral Ward 1 compared to Wards 2 and 6
124 (Figure 1B) and the proportion of SARS-CoV-2 PCR tests with positive results was approximately five-
125 fold higher (Supplementary Figure S3A). Both of these proxy measures of epidemic intensity (θ_j^{PCR} and

126 PCR test positivity) are limited by potential biases, some of which are likely to still be present even if
127 PCR testing rates are relatively equal by ward [16]. Nonetheless, these measures suggest substantial
128 heterogeneity in the underlying epidemic intensity, with an apparent higher rate of previously infected
129 individuals (as a proportion of the population) in Somerville Wards 1 and 4. Thus, we observed that the
130 venue location chosen for this study, and its associated survey participant catchment distribution,
131 resulted in relative under-sampling of wards with expected higher seropositivity and oversampling of
132 those with lower expected seropositivity (Figure 1C).

133 *GPS-estimated visitor catchment distributions*

134 Recognizing that survey participant catchment distributions can be poorly matched to the underlying
135 geographic distribution of seropositivity, we explored the use of GPS-estimated foot traffic data as a tool
136 for evaluating actual or candidate locations for venue-based sampling. We evaluated correlations
137 between the observed participant catchment distribution for the actual Somerville study location, GPS-
138 estimated visitor catchment distributions for this site and a hypothetical alternative site in Somerville
139 Ward 1, and the cumulative incidence of PCR-confirmed infections by ward. The participant catchment
140 distribution at the actual study site (P_j^{direct}) closely matched its corresponding GPS-estimated visitor
141 catchment distribution, V_j^{site} (Pearson's $r = 0.90$, $p = 0.0131$, Figure S4). However, the GPS-estimated
142 visitor catchment distribution for the actual study site was poorly correlated with the cumulative
143 incidence of PCR-confirmed infections ($r = -0.11$, $p = 0.55$, Figure 2C).

144 We evaluated whether choosing an alternative study site could improve the correlation between sample
145 allocation and cumulative incidence of PCR-confirmed infections by ward, with the goal of reducing
146 uncertainty in estimated seroprevalence. We observed that the GPS-estimated visitor catchment
147 distribution at the alternative site (V_j^{alt}) is strongly correlated with ward-level cumulative incidence of
148 PCR-confirmed infections ($r = 0.93$, $p = 0.0072$, Figure 2D). If subgroup sizes are known and differences
149 in subgroup-level seropositivity can be inferred or assumed, allocating more samples to larger subgroups

150 and those with higher expected seropositivity will improve precision for weighted population-level
151 seroprevalence estimates (Figure S5 and Equation 2) [17]. This suggests that, if the GPS-estimated
152 visitor catchment distribution reliably predicts the survey participant catchment at the alternative study
153 site, this location would yield improved sample allocation across geographic subgroups that more
154 closely approximates optimal sample allocation.

155 *Venue location and uncertainty in SARS-CoV-2 seroprevalence estimates*

156 We used numerical simulation to quantify uncertainty in estimated SARS-CoV-2 seroprevalence under
157 different survey participant catchment distributions. Specifically, we constructed a simple synthetic
158 population with seven geographic subgroups, corresponding to the seven electoral wards in Somerville,
159 each divided into three age-based subgroups. We specified the size of each subgroup using local census
160 data [15] and specified the true underlying seropositivity for each age-location subgroup by assuming
161 these values are proportional to the observed cumulative incidence of PCR-confirmed infections for
162 each subgroup. Using this model, we compared three sample allocation scenarios: (1) optimal allocation,
163 in which the number of individuals sampled from each age-location subgroup is specified to optimally
164 reduce uncertainty in the resulting seroprevalence estimates (per Equation 2 in the Supplementary
165 Information); (2) allocation according to the observed survey participant catchment distribution for the
166 actual study site; (3) allocation according to the GPS-estimated visitor distribution at the hypothetical
167 alternative study site. (Additional details on model specification and sensitivity testing for model
168 parameters are available in the Supplementary Information.)

169 We observed 1.5- to 2-fold higher uncertainty when sampling effort was allocated according to the
170 participant catchment distribution at the study site compared to the alternative site or optimal allocation
171 (Figure 3). This observation suggests that choice of recruitment location can result in suboptimal sample
172 allocation and higher uncertainty.

173 *Bias due to unappreciated heterogeneity in seropositivity across geographic subgroups*

174 Biased seroprevalence estimates can result if geographic heterogeneity in sample allocation results in
175 substantial over- or under-sampling in locations with higher (or lower) seropositivity, and if procedures
176 for generating weighted prevalence estimates do not appropriately account for geographic heterogeneity
177 in underlying seropositivity. We compared estimated seroprevalence versus true seroprevalence for
178 numerical simulations in which the final seroprevalence estimates were weighted by the sampling
179 probability for each age-location group or by the sampling probability of each age subgroup alone
180 (Figure 4). The first weighting procedure accounts for heterogeneity across age and location subgroups,
181 whereas the second procedure accounts only for heterogeneity across age subgroups. Using the second
182 procedure resulted in over- or under-estimation of seroprevalence, depending on whether sample
183 allocation enriches for participants from areas with high or low underlying seropositivity, respectively.

184 **Discussion**

185 Convenience sampling, despite its inherent limitations, may have continued utility in the public health
186 response to the COVID-19 pandemic and for infectious disease outbreaks generally. Cost and logistical
187 considerations may limit the feasibility of randomized structured sampling, particularly in resource-
188 constrained contexts or in situations where census data, population rosters, or household mapping data
189 are unavailable or unreliable. Certain forms of convenience sampling are better suited for reaching
190 important subgroups compared to structured approaches. Lower-wage or frontline workers who are at
191 higher risk of SARS-CoV-2 exposure [18-20], including undocumented workers [20], may be less likely
192 to participate if recruited using conventional survey outreach methods (e.g., mail or phone contact) due
193 to constraints on their time [21-23] and lack of incentives [21]. Convenience sampling at highly visited
194 community locations such as essential businesses may be an attractive alternative to structured sampling
195 in this important population, similar to sampling approaches developed to study so-called “hidden
196 populations” [24].

197 Geographic heterogeneity in SARS-CoV-2 epidemic intensity within cities has been a repeatedly
198 observed feature of the pandemic [1–3]. This phenomenon poses unique challenges for seroprevalence
199 studies that employ venue-based and other convenience sampling strategies, in which sample allocation
200 across subgroups cannot be pre-specified and is non-uniform across geographic space.

201 Examining data from our seroprevalence study in Somerville, MA, we observed that venue-based
202 sampling resulted in substantial undersampling of areas where proxy measures (cumulative incidence of
203 PCR-confirmed infections and PCR test positivity rates) suggest higher epidemic intensity. This
204 mismatch between the survey participant catchment distribution and the geographic distribution of
205 seropositive individuals can result in suboptimal sample allocation and higher variance in resultant
206 seroprevalence estimates.

207 Study locations can have widely divergent participant catchment distributions. These distributions can
208 be heavily enriched for participants living in the immediate vicinity of the study location. This limitation
209 has important implications for bias and uncertainty of resulting seroprevalence estimates (as detailed
210 here) and raises questions about potential undersampling or exclusion of important subgroups in venue-
211 based studies. Multiple studies have identified geographic location as a strong surrogate for multiple risk
212 factors associated with severe infection, hospitalization, and/or death due to COVID-19 [25] and
213 undersampling in neighborhoods where these risk factors co-localize together can compromise the
214 reliability and interpretability of seroprevalence estimates. Recruiting participants directly from such
215 communities, where rates of COVID-19 related hospitalization and deaths are often higher, has yielded
216 seroprevalence estimates that are substantially higher than city-level or state-level estimates [4,9].

217 Notably, the areas that were most undersampled in our study strongly overlap neighborhoods with lower
218 socioeconomic status, larger proportions of non-white residents, lower proportions of English-speaking
219 households (Figure 1A, Supplementary Figure S3C).

220 Our work has three practical findings that are applicable to the design, implementation, and
221 interpretation of convenience-based seroprevalence studies.

222 (1) Uncertainty in population-level seroprevalence estimates is minimized when sample allocation is
223 proportional to the size and underlying seropositivity of individual subgroups in the population
224 (Equation 2 and Figure S5). The practical application of this finding may be limited because of
225 challenges in reliably ascertaining differences in the underlying seroprevalence between subgroups *a*
226 *priori*. For example, this may arise when access to diagnostic testing for acute infections is limited or
227 disparate across subgroups. However, even in this situation, allocating sampling effort proportional to
228 subgroup sizes alone can substantially reduce uncertainty (Figure S1). Applying this finding may also be
229 difficult because an inherent feature of venue-based sampling is that allocation of sampling effort is not
230 pre-specified, but instead results from a stochastic process that depends on the location of the study
231 venue. Thus, although optimal sample allocation is likely not achievable via venue-based sampling,
232 careful selection of venue locations, with the objective of enriching for participants from geographic
233 subgroups with larger populations and/or higher expected seroprevalence, can at least improve sample
234 allocation and help reduce uncertainty.

235 (2) GPS-estimated foot traffic can inform the selection of venue-based recruitment locations. The GPS-
236 estimated visitor catchment distribution at our study location correlated closely with the survey
237 participant catchment distribution. Validation against other data sources that directly measure the
238 geographic distributions of visitors to locations of interest (for example, aggregated geographic and
239 registration data from COVID-19 mobile testing programs) can help further evaluate this potentially
240 important data source.

241 (3) Convenience sampling can produce biased seroprevalence estimates if geographic heterogeneity in
242 underlying subgroup-level seropositivity is not properly accounted for (Figure 5). Consistent with prior
243 studies [2, 3], we observed that COVID-19 incidence can vary widely, even over a relatively small
244 geographic area. To avoid this problem, studies that employ convenience sampling should collect
245 geographic data on participants' home locations that is granular enough to capture potential geographic
246 heterogeneity in seroprevalence within the study area. This information, combined with data on the

247 catchment distributions of individual recruitment sites (including GPS-estimated foot traffic data, as
248 examined here), can be used to quantify what would otherwise be an unmeasured source of bias in
249 resulting seroprevalence estimates.

250 Multiple considerations are important for contextualizing our findings and the recommendations above.
251 These include participation bias that may result in exclusion of individuals with disabilities or others
252 who are less likely to leave their homes. Methods designed to account for participation bias, including
253 those developed for use with time-location sampling [26], may be applicable here. Likewise, collecting
254 information on non-respondents in venue-based sampling—for example, brief demographic surveys
255 collected before recruitment for serological testing—can help measure and account for potential sources
256 of participation bias.

257 The GPS-estimated foot traffic data used in our study have several important limitations. Identification
258 of visitors and their home locations in this data may be biased by differences in mobile device usage
259 between demographic groups, potentially under-sampling visitors from important populations or
260 oversampling others. In addition, this data can be sparse and thus more subject to stochastic variation
261 when only small numbers of users are captured in either the point of interest or home location. Different
262 forms of participation bias (described above) are expected to skew the geographic distribution of study
263 participants away from GPS-estimated catchment distributions, potentially making it difficult to use this
264 data source in real-world public health practice.

265 Lastly, we make several simplifying assumptions in our numerical model. The numerical model assumes
266 that the true seropositivity in each age-location subgroup is proportional to its observed cumulative
267 incidence of PCR-confirmed SARS-CoV-2 infections (per local public health data from Somerville,
268 MA). However, wards with higher PCR positivity rates (an indicator of greater epidemic intensity) have
269 relatively the same rates of overall PCR testing per capita (Supplementary Figure S3B), indicating that
270 there were gaps in testing effort in areas of Somerville with more incident infections overall [16]. The
271 assumed true underlying seroprevalence of each age-location group, which is specified using the

272 observed cumulative incidence of PCR-confirmed infection and does not account for the testing gap
273 described above, are less dispersed across age-location groups than what would be expected if PCR
274 testing effort better matched epidemic intensity by ward (i.e., greater PCR testing effort in heavily
275 impacted areas would likely reveal even larger differences in cumulative incidence between wards). This
276 misspecification, and resultant smaller dispersion in assumed true cumulative incidence by ward, is
277 expected to result in more conservative values for the uncertainty in estimated population-level
278 seroprevalence; otherwise, this limitation is not expected to change our primary findings from the
279 numerical model.

280 In summary, we have examined how geographic heterogeneity in sample allocation, combined with
281 underlying heterogeneity in geographic distribution of seropositive individuals, can influence
282 seroprevalence estimates derived from venue-based sampling. Our findings are relevant to studies
283 employing venue-based recruitment and are also applicable to other kinds of convenience sampling, for
284 example, studies using a hospital's discarded blood specimens from patients drawn from the hospital's
285 geographic catchment area.

286 **References**

- 287 1. Feehan AK, Fort D, Garcia-Diaz J, et al. Seroprevalence of SARS-CoV-2 and Infection Fatality
288 Ratio, Orleans and Jefferson Parishes, Louisiana, USA, May 2020. *Emerging Infectious Diseases* **2020**;
289 26(11).
- 290 2. Kim SJ, Bostwick W. Social Vulnerability and Racial Inequality in COVID-19 Deaths in
291 Chicago. *Health Education & Behavior* **2020**; 47(4):509–513.
- 292 3. Kissler S, Kishore N, Prabhu M, et al. Reductions in commuting mobility predict geographic
293 differences in SARS-CoV-2 prevalence in New York City. *Nature Communications* **2020**; 11(1):4674.

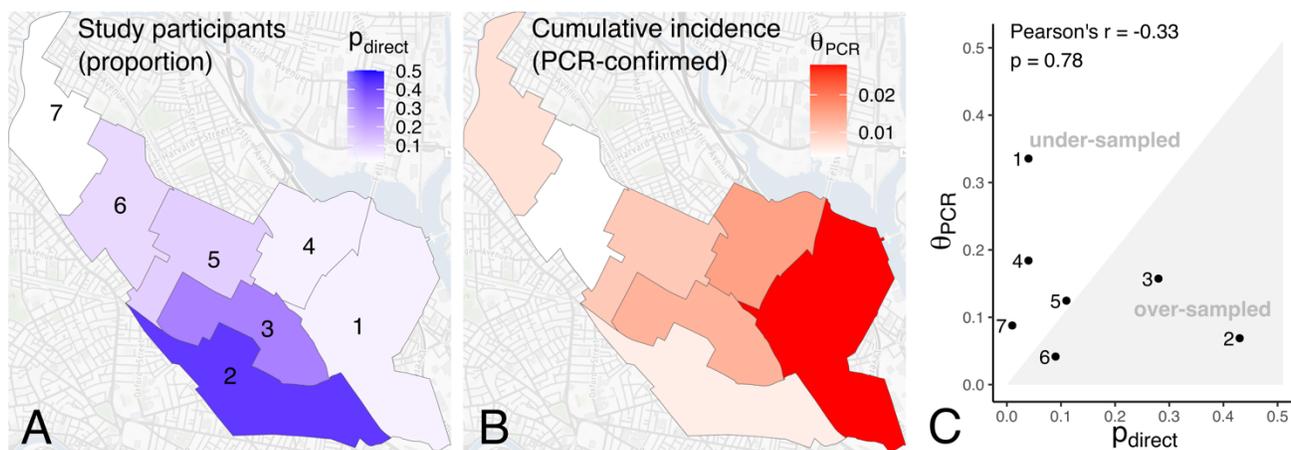
- 294 4. Rosenberg ES, Tesoriero JM, Rosenthal EM, et al. Cumulative incidence and diagnosis of
295 SARS-CoV-2 infection in New York. *Annals of Epidemiology*. **2020**; 48:23–29
- 296 5. Havers FP, Reed C, Lim T, et al. Seroprevalence of Antibodies to SARS-CoV-2 in 10 Sites in
297 the United States, March 23-May 12, 2020. *JAMA Internal Medicine* **2020**; 180(12):1576–1586
- 298 6. Bendavid E, Mulaney B, Sood N, et al. COVID-19 Antibody Seroprevalence in Santa Clara
299 County, California. *International Journal of Epidemiology* **2021**; 50(2):410-419.
- 300 7. Shook-Sa BE, Boyce RM, Aiello AE. Estimation Without Representation: Early Severe Acute
301 Respiratory Syndrome Coronavirus 2 Seroprevalence Studies and the Path Forward. *J Inf Dis*. **2020**;
302 222(7):1086–1089.
- 303 8. Kelly H. A random cluster survey and a convenience sample give comparable estimates of
304 immunity to vaccine preventable diseases in children of school age in Victoria, Australia. *Vaccine*.
305 **2002**; 20(25-26):3130–3136.
- 306 9. Naranbhai V, Chang CC, Beltran WFG, et al. High seroprevalence of anti-SARS-CoV-2
307 antibodies in Chelsea, Massachusetts. *J Inf Dis* **2020**; 222(12):1955–1959.
- 308 10. Chang, S., Pierson, E., Koh, P.W. et al. Mobility network models of COVID-19 explain
309 inequities and inform reopening. *Nature* **2021**; 589: 82–87.
- 310 11. Ma, KC and Lipsitch M. Big data and simple models used to track the spread of COVID-19 in
311 cities. *Nature* **2021**; 589: 26-28.
- 312 12. Hobbs CV, Drobeniuc J, Kittle T, et al. Estimated SARS-CoV-2 seroprevalence among persons
313 aged <18 years - Mississippi, May-September 2020. *MMWR Morbidity and Mortality Weekly Report*.
314 **2021**; 70(9):312–315.

- 315 13. Sutton M, Cieslak P, Linder M. Notes from the field: Seroprevalence estimates of SARS-CoV-2
316 infection in convenience sample - Oregon, May 11-June 15, 2020. *MMWR Morbidity and Mortality*
317 *Weekly Report.* **2020**; 69(32):1100–1101.
- 318 14. Massachusetts State Department of Health. COVID-19 Response Reporting. Available from:
319 <https://www.mass.gov/info-details/covid-19-response-reporting>
- 320 15. City of Somerville, Massachusetts and Cambridge Health Alliance. The Wellbeing of Somerville
321 Report. Available at: [https://www.somervillema.gov/sites/default/files/wellbeing-of-somerville-report-](https://www.somervillema.gov/sites/default/files/wellbeing-of-somerville-report-2017.pdf)
322 [2017.pdf](https://www.somervillema.gov/sites/default/files/wellbeing-of-somerville-report-2017.pdf)
- 323 16. Dryden-Peterson S, Velásquez GE, Stopka TJ, Davey S, Lockman S, Ojikutu B. Disparities in
324 SARS-CoV-2 Testing in Massachusetts During the COVID-19 Pandemic. *JAMA Network Open* **2020**;
325 4(2):e2037067
- 326 17. Larremore DB, Fosdick BK, Bubar KM, et al. Estimating SARS-CoV-2 seroprevalence and
327 epidemiological parameters with uncertainty from serological surveys. *Elife* **2021**; 10:e64206
- 328 18. Hawkins D. Social Determinants of COVID-19 in Massachusetts, United States: An Ecological
329 Study. *Journal of Preventive Medicine and Public Health.* Korean Society for Preventive Medicine;
330 **2020**; 53(4):220–227.
- 331 19. Baker MG, Peckham TK, Seixas NS. Estimating the burden of United States workers exposed to
332 infection or disease: A key factor in containing risk of COVID-19 infection. *PLOS ONE* **2020**;
333 15(4):e0232452.
- 334 20. Feehan AK, Velasco C, Fort D, et al. Racial and workplace disparities in seroprevalence of
335 SARS-CoV-2 in Baton Rouge, Louisiana, July 15-31, 2020. *Emerging Infectious Diseases* **2021**; 27(1):
336 314-317

- 337 21. Hernández MG, Nguyen J, Casanova S, Suárez-Orozco C, Saetermoe CL. Doing no harm and
338 getting it right: Guidelines for ethical research with immigrant communities. *New Directions for Child*
339 *and Adolescent Development*. **2013**; 2013(141):43–60.
- 340 22. Corbie-Smith DM. Minority recruitment and participation in health research. *North Carolina*
341 *Medical Journal*. **2004**; 65(6):385-7
- 342 23. Keyzer JF, Melnikow J, Kuppermann M, et al. Recruitment strategies for minority participation:
343 challenges and cost lessons from the POWER interview. *Ethn Dis* **2005**; 15(3):395–406.
- 344 24. Muhib FB, Lin LS, Stueve A, et al. A Venue-Based Method for Sampling Hard-to-Reach
345 Populations. *Public Health Reports* **2001**; 116:216–222.
- 346 25. Maroko AR, Nash D, Pavilonis BT. COVID-19 and inequity: A comparative spatial analysis of
347 New York City and Chicago hot spots. *Journal of Urban Health* **2020**; 97(4):461–470.
- 348 26. Leon L, Jauffret-Roustide M, Le Strat Y. Design-based inference in time-location sampling.
349 *Biostatistics* **2015**; 16(3):565–579.

350

351 **Figures**

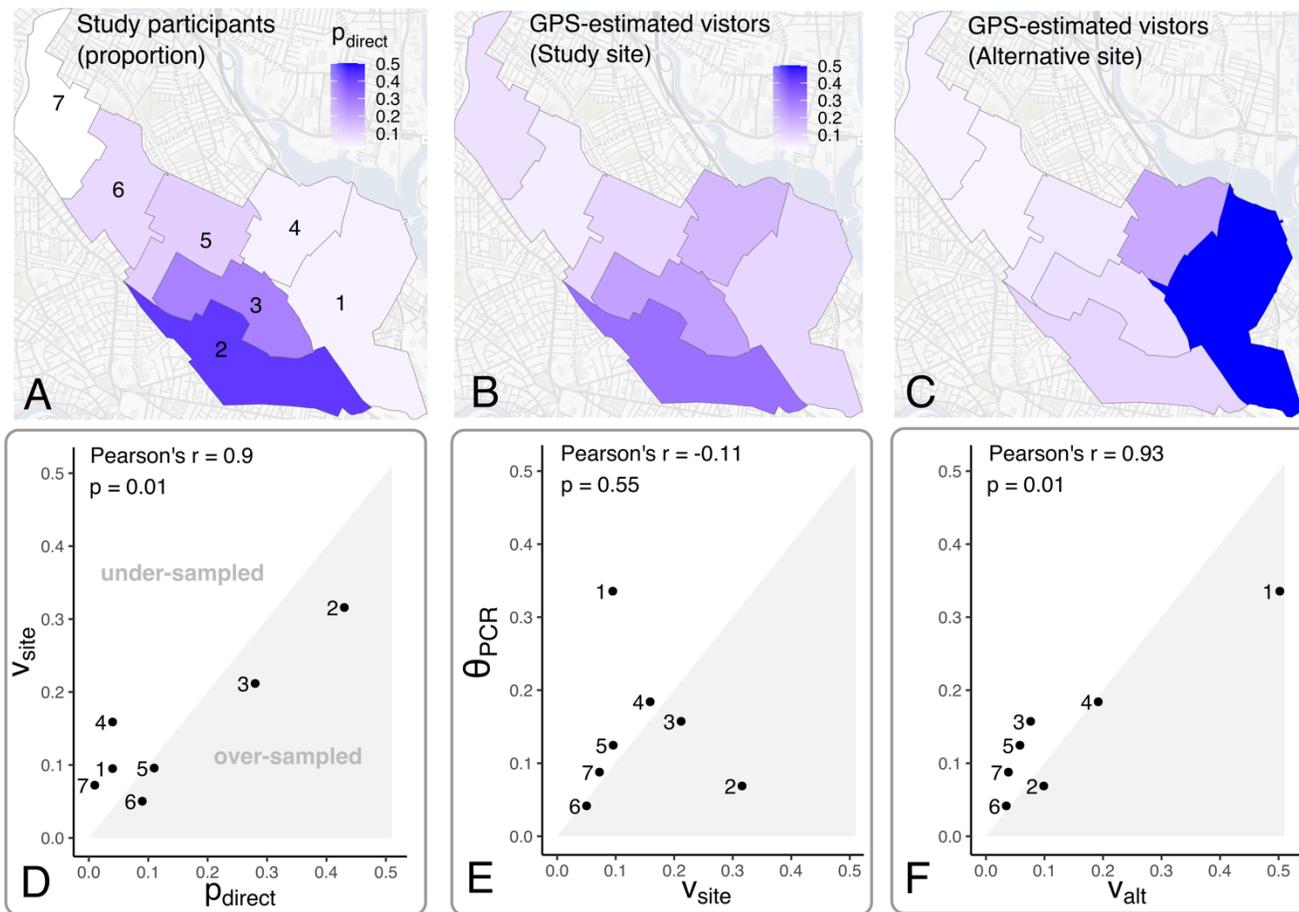


352

353 **Figure 1. Sample allocation and geographic heterogeneity in proxy measures of epidemic**

354 **intensity.** (A) Survey participant catchment distribution. Wards are shaded by p_j^{direct} , the proportion of
355 all directly recruited participants from each of Somerville Wards 1-7; (B) Cumulative incidence of prior
356 PCR-confirmed SARS-CoV-2 infections by Ward as of June 8th, 2020 (θ_j^{PCR}); (C) Correlation between
357 p_j^{direct} and θ_j^{PCR} . Significance of the correlation is calculated via permutation testing, as described in the
358 Supplementary Material.

359



360

361 **Figure 2. GPS-estimated visitor catchment distributions for actual and hypothetical alternative**

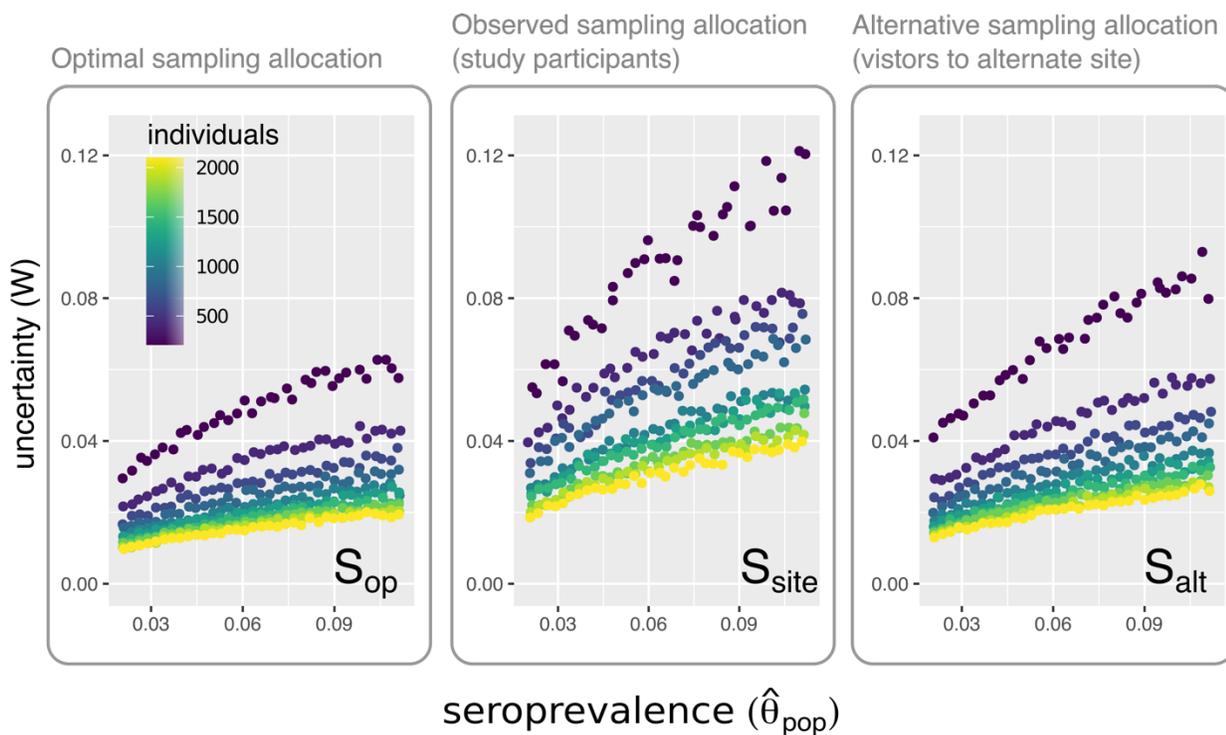
362 **study sites.** GPS-estimated visitor catchment distributions for (A) the actual study location V_j^{site} or (B)

363 a hypothetical alternative study site in Somerville Ward 1 V_j^{alt} . Correlations between V_j^{site} , V_j^{alt} , and the

364 cumulative incidence of PCR-confirmed SARS-CoV-2 infection θ_j^{PCR} are shown in (C) and (D).

365

366



367

368 **Figure 3. Uncertainty in estimated SARS-CoV-2 seroprevalence obtained using different**

369 **sample allocation strategies.** The uncertainty (W , the width of the 95th percentile interval for

370 10000 estimated seroprevalence values) versus mean estimated seroprevalence for different values

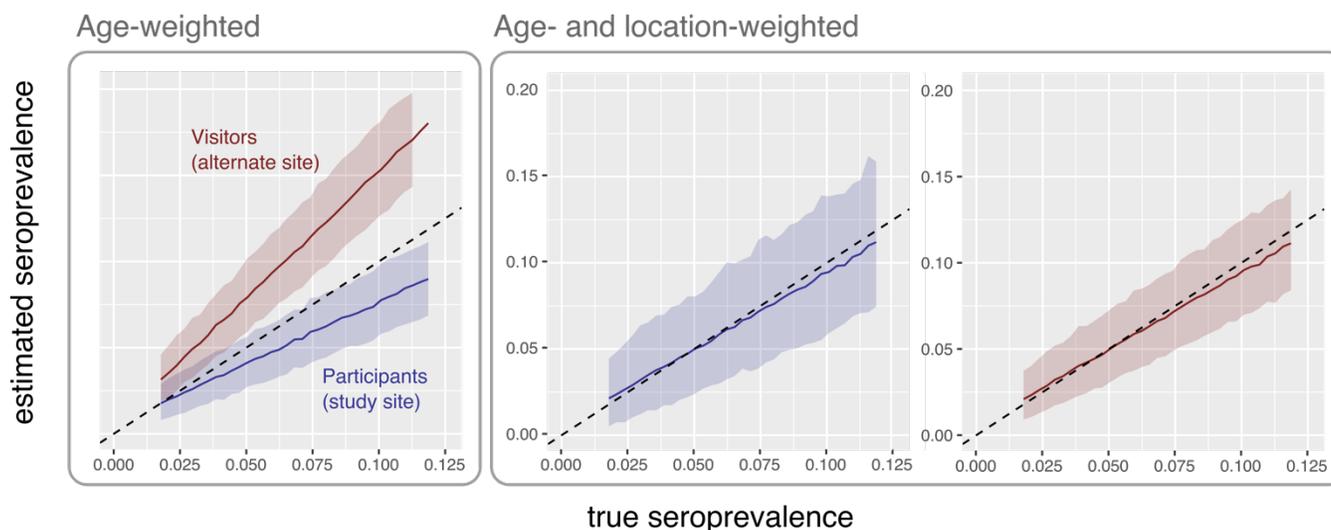
371 of n (the total number of individuals sampled) when individuals are sampled using (1) in the left

372 panel, the optimal sample allocation according to Equation 2 in the Supplementary Material (S_{op});

373 (2), in the center panel, the sampling distribution of participants in the Somerville seroprevalence

374 survey (S_{site}); (3) the sampling distribution at the proposed alternative study site (S_{alt}).

375



376

377 **Figure 4. Bias in estimated seroprevalence by sampling strategy and weighting procedures.** Left:

378 Estimated seroprevalence, weighted only by age subgroups. Right: Estimates weighted by age-location

379 subgroups are shown in the right panel. Blue: Sample allocation specified by the observed participant

380 distribution catchment distribution in the Somerville study. Red: Sample allocation specified by the

381 catchment distribution of GPS-estimated visitors to the proposed alternate study site. Dotted line

382 indicates where estimated equals true seroprevalence.

383

384 **Funding**

385 This work was supported by the Andrew and Corey Morris-Singer Foundation, National Cancer Institute

386 at the National Institutes of Health [U01CA261277] and the National Institute of Allergy and Infectious

387 Diseases at the National Institutes of Health [T32AI007061 to TSB]. This project has been funded in

388 part by contract 200-2016-91779 with the Centers for Disease Control and Prevention. The findings,

389 conclusions, and views expressed in this presentation are those of the author(s) and do not necessarily

390 represent the official position of the Centers for Disease Control and Prevention (CDC).

391