

1 Bayesian adaptive design for citizen science data
2 collection: Exploring tensions between data and design

3 Max Savery¹ and Stijn Luca¹

4 ¹Ghent University, Department of Data Analysis and Mathematical
5 Modelling
6 Coupure links 653, B-9000 Ghent, Belgium

7 **1 Abstract**

8 1. Bayesian adaptive design can be applied in spatial settings where future survey locations
9 need to be selected based on already available data. An important use case of adaptive
10 design is the recommendation of locations for opportunistic, citizen science collection of
11 species observation data, where some areas are already overrepresented and others are
12 severely undersampled.

13 2. This work conducts an extensive comparative study of adaptive Bayesian design
14 approaches, specifically for recommending survey locations for citizen science data
15 collection. We explore compromises between design-based, model-based, and
16 exchange-based approaches in order to better understand the statistical implications of
17 using adaptive design for citizen science projects.

18 3. To evaluate the adaptive design approaches effectively, we work in a simulated
19 environment, allowing us to compare the performance between fourteen different design

20 strategies. For each method, we vary design size and the initial data conditions, and in a
21 separate analysis show the effect of weighting the site recommendations by varying degrees
22 of citizen scientist preference. Additionally, we conduct a sensitivity experiment on the
23 approximation used to integrate out the uncertainty associated with the future data
24 collection.

25 4. Our results highlight the tension between the data and the design methodology: We
26 show that the model- and exchange-based approaches perform better when more prior data
27 is available, particularly when estimation or prediction accuracy is preferred. In contrast,
28 design-based methods are more stable when prior information is limited. However, it
29 remains a challenge to balance all design priorities and data scenarios in a single approach.
30 Additionally, we show that including citizen scientist preference in the utility function can
31 impact the design in surprising ways. These outcomes are not uniform and depend on the
32 method used to optimize the design, the metric of interest for optimization, as well as the
33 amount of prior data and the desired size of the future survey. Taken together, our results
34 show that effective survey design must be carefully matched to both the available baseline
35 information and the study objectives.

36 **Keywords**

37 bayesian adaptive design, adaptive sampling, spatial design, opportunistic data collection,
38 citizen science, biodiversity monitoring, species distribution modelling

2 Introduction

Traditional surveys for estimating species distributions are designed and carried out according to an imposed structure, ensuring that a statistical model can be fit without bias in the parameter estimation. However, unstructured citizen science data has become a useful and alternative source of species observation data [Zipkin et al., 2021]. For example, citizen science presence-only (PO) data has been shown to be a valuable supplement to presence-absence (PA) survey data [Dorazio, 2014, Koshkina et al., 2017, Simmonds et al., 2020] when both data sources are combined in an integrated model. Though such citizen science data is ubiquitous and has demonstrated benefit for ecological modelling, it exhibits highly opportunistic data collection patterns [Johnston et al., 2023]. It is therefore desirable to direct citizen scientists to locations that maximize information complementary to the data already collected.

To recommend locations, the framework of Bayesian adaptive design can be applied to survey design for species monitoring. Generally, Bayesian design is used to find a configuration of design points that optimize a utility or loss function of interest. Adaptive design can be considered a useful extension when previously collected data is available [Mateu and Muller, 2013]. In the species monitoring setting, this is framed as the construction of a survey conditioned on data previously collected at a limited set of locations. The design itself is characterized by the number of sites to visit, the locations of the sites, and the number of visits to make at each site [Reich, 2020]. However, there is a limited body of research and no standard guidance regarding design strategies for citizen scientists—Bayesian, adaptive, or otherwise. Based on the few works that do consider survey design, there appears to be substantial advantages to including adaptive design

62 methodology in citizen science data collection. Callaghan et al. [2019] use statistical
63 leverage to assign a marginal value to each sampling event, in order to predict leverage
64 values for future sampling events, and they highlight the importance of sampling both
65 hotspots and unsampled areas, depending on desired design resolution. Mondain-Monval
66 et al. [2024] compare six design- and model-based approaches for adaptive design. They
67 target uncertainty and rare species in their designs and consider varying levels of citizen
68 science uptake of the recommended locations. They find that even a small proportion of
69 adaptively sampled data improved species distribution model performance. Flint et al.
70 [2024], using citizen science roadkill as a case study, study the utility of additional
71 sampling by comparing relative intensities between models with and without predictions
72 made at an additional design point. Importantly, they incorporate organizational value and
73 financial cost of additional sampling into their utility function. They highlight that their
74 framework can avoid resource waste and select high-utility locations that are particularly
75 valuable for future sampling. While these works provide excellent practical approaches to
76 the design of adaptive surveys for citizen scientists and consider many factors important for
77 citizen science surveying, they do not, however, embed their methods within the broader
78 framework of Bayesian design. Furthermore, they do not offer multifaceted comparisons of
79 algorithms for adaptive design, nor do they provide a general framework for including
80 citizen scientist preference in the design optimization.

81 Beyond citizen scientist data collection, there are a few other notable ecological
82 applications of Bayesian adaptive design. Williams et al. [2018] use the Bayesian adaptive
83 design framework for creating designs that measure dynamic populations. In their work,
84 future data is simulated from the posterior predictive distribution of prior data, in order to

85 minimize abundance uncertainty and determine the optimal next site to visit. Similarly,
86 Leach et al. [2022] formalize the recursive nature of Bayesian adaptive design for ecological
87 applications, specifically when using the Prior-Proposal Recursive Bayes approach of
88 Hooten et al. [2021]. Thilan et al. [2023] and Thilan et al. [2024] evaluate the Bayesian
89 adaptive design framework in its use for tracking changes and trends in ecological
90 monitoring, using a Laplace approximation to estimate the posterior distribution and
91 assess the change in information when adding new sites to the already collected data with
92 Kullback-Leibler (KL) divergence.

93 When planning and managing citizen science projects, not only are the statistical
94 properties of the design important, it is also essential to consider citizen scientist
95 behaviour. Callaghan et al. [2023] study the extent to which citizen scientists are willing to
96 follow such adaptive sampling schemes and show that "behavioural nudges" do in fact
97 incentivize citizen scientist to survey recommended locations. Pocock et al. [2023] study
98 the effect of different types of citizen scientist recorder behaviour on the estimates of
99 species trend and occupancy across different butterfly species in the UK. August et al.
100 [2020] characterize the behaviour of citizen scientists using a variety of derived metrics
101 based on aspects of the data collection in space, time, and content. These scores are then
102 clustered and the principal components analysed in order to describe the behaviour via
103 continuous axes. In addition, Callaghan and Gawlik [2015] compare estimates of species
104 abundance derived from eBird data and from structured survey data, finding citizen science
105 surveys can be used as an effective substitute for the structured surveys. Callaghan et al.
106 [2022] focus on sampling effort, highlighting the relationship between the amount of
107 collected citizen science data, sample completeness, and species diversity. Thompson et al.

108 [2023] survey participants of a citizen science project and find that some people may be
109 willing to change their surveying behaviour if they knew the change would lead to data
110 more valuable for the project. However, while the works mentioned here do address citizen
111 scientist behaviour and the impact of citizen science data collection on downstream
112 modelling outcomes, they do not fully address the statistical issues related to adaptive
113 design for citizen scientists. Henrys et al. [2024] review at a broad level the issues with
114 creating adaptive designs for citizen scientists, suggesting that clearer guidance regarding
115 the statistical methodology and additional examples of its application are required for the
116 uptake of adaptive design methodology.

117 The literature described above indicates that there are certainly advantages of using the
118 framework of Bayesian adaptive design to recommend site locations for citizen science data
119 collection. It is also evident that while there are a number of potential approaches for
120 Bayesian adaptive designs for citizen scientists, the advantages, disadvantages, and
121 practicalities of implementing Bayesian design algorithms are not well-studied. This leaves
122 us with the question as to what works best for citizen science applications, what
123 refinements are needed in the Bayesian design framework, and how can we include the
124 most information from both a statistical perspective and a practical one when generating
125 designs and recommending potential survey locations that can be provided either directly
126 to an interested citizen scientist or via an organization tasked with developing a citizen
127 science survey project. Building on the adaptive design work of Callaghan et al. [2019],
128 Flint et al. [2024], and Mondain-Monval et al. [2024], and considering the gaps in the
129 literature discussed by Henrys et al. [2024], in this study we explore the associated
130 statistical advantages and complications that come with using Bayesian adaptive design for

131 citizen science applications. We expand specifically upon Mondain-Monval et al. [2024] by
132 conducting a more extensive comparison of design- and model-based methods, as well as
133 the exchange family of algorithms, and we consider the practical challenges of
134 implementing model-based and exchanged-based approaches in a Bayesian design
135 framework. An additional novel contribution of our work is to show the effect of
136 considering citizen scientist accessibility or preference on the utility functions used with the
137 design optimization. Importantly, we also perform a sensitivity analysis on the effect of
138 drawing future data from the posterior predictive distribution of the model fit to the initial
139 observations, a necessary step in Bayesian adaptive design but one that is rarely assessed.
140 Finally, to conclude our comparative study, we provide a list of five high-level
141 recommendations for the adaptive design of citizen science surveys. These
142 recommendations take into account the the compromise between the computation efficiency
143 of design-based approaches and the flexible optimization of the more computationally
144 intensive Bayesian model- or exchange based algorithms, as well as the effect of
145 incorporating citizen scientist preference in the optimization of the designs themselves.
146 By putting these aspects of design and citizen science data collection together in a single
147 work, we seek to make Bayesian adaptive design more transparent and accessible to
148 organizations with interest in designing surveys for citizen scientists.

149 **3 Methods**

150 In the this section we introduce the framework of Bayesian adaptive design and describe
151 the design approaches we use for creating Bayesian adaptive designs for citizen scientists.

3.1 Citizen science data

We refer to a PO (presence-only) dataset as a set of points $\mathbf{y}_0 = \{s_1, s_2, \dots, s_n\}$ collected in region \mathcal{D} , with n total observations in the dataset. The region D is divided into c cells of equal area $|A_i| = |D|/c$, where A_i refers to each cell or site. Following Mondain-Monval et al. [2024] we use a virtual ecologist approach [Zurell et al., 2010] to study the effect of different adaptive design methodologies on citizen science data collection and species distribution model (SDM) estimation. We use the entire domain of the country of Belgium as a case-study for survey design, simulating species observations via covariates derived from remote-sensing data within the country. We briefly summarize the initial data generation workflow as follows. The data generation and model estimation details are more fully described in the Supplementary material.

1. We use a Log-Gaussian Cox Process (LGCP) as the true data-generating SDM. This data-generating model is used to generate the initial dataset \mathbf{y}_0 , which we treat as the initial, fixed set of citizen science observations. In practice, the \mathbf{y}_0 dataset will be the previously collected data upon which the adaptive design is built.
2. We then fit a working model to the simulated \mathbf{y}_0 data. At this step, we again use the LGCP, and we induce a misspecification of the bias function in order to mimic situations where the data generating process is not fully known or model specification may be difficult.
3. Using the posterior from the working model, we simulate R posterior predictive datasets, referred to as $\tilde{\mathbf{y}}_1^{(r)}$. These R datasets represent possible future data outcomes and are used within the Bayesian design framework to approximate the

174 expectation of the utility function.

175 In practical applications of Bayesian design, it will be necessary to use the posterior
176 resulting from \mathbf{y}_0 to generate the $\tilde{\mathbf{y}}_1$ datasets, and the working procedure that we use here
177 reflects that reality. This is explained in more detail in the next section.

178 **3.2 Bayesian adaptive design**

179 The purpose of Bayesian adaptive design is to find the configuration of design points \mathbf{d}^*
180 within region \mathcal{D} that maximizes the expected utility or minimizes the expected loss taken
181 over the joint distribution $p(\boldsymbol{\theta}, \mathbf{y}_1 | \mathbf{y}_0, \mathbf{d})$ of parameters $\boldsymbol{\theta}$ and future data \mathbf{y}_1 , conditional on
182 the prior data \mathbf{y}_0 . In this work we primarily refer to maximization of the expected utility
183 function $E[U(\mathbf{d}) | \mathbf{y}_0]$:

$$\mathbf{d}^* = \operatorname{argmax}_{\mathbf{d} \in \mathcal{D}} E[U(\mathbf{d}) | \mathbf{y}_0]. \quad (1)$$

184 The Bayesian adaptive design framework accounts for the uncertainty associated with both
185 future data collection and parameter estimation associated with the joint distribution.
186 However, a primary challenge is the approximation of the expectation of the utility for
187 every proposed design. The expected utility is defined as

$$E[U(\mathbf{d}) | \mathbf{y}_0] = \int_{\mathbf{y}_1} \int_{\boldsymbol{\theta}} u(\mathbf{y}_1, \boldsymbol{\theta}, \mathbf{d}) p(\boldsymbol{\theta} | \mathbf{y}_1, \mathbf{y}_0, \mathbf{d}) p(\mathbf{y}_1 | \mathbf{y}_0, \mathbf{d}) d\boldsymbol{\theta} d\mathbf{y}_1, \quad (2)$$

188 where $p(\boldsymbol{\theta} | \mathbf{y}_1, \mathbf{y}_0, \mathbf{d})$ is the posterior of model parameters, and $p(\mathbf{y}_1 | \mathbf{y}_0, \mathbf{d})$ is the posterior
189 predictive distribution (PPD) induced by the historical data \mathbf{y}_0 , where
190 $p(\mathbf{y}_1 | \mathbf{y}_0, \mathbf{d}) = \int_{\boldsymbol{\theta}} p(\mathbf{y}_1 | \boldsymbol{\theta}, \mathbf{d}) p(\boldsymbol{\theta} | \mathbf{y}_0) d\boldsymbol{\theta}$. This distinguishes adaptive design from non-adaptive

191 design where the prior predictive distribution $p(\mathbf{y}|\mathbf{d}) = \int_{\boldsymbol{\theta}} p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{d})p(\boldsymbol{\theta})d\boldsymbol{\theta}$ is used instead,
 192 with no distinction between \mathbf{y}_1 and \mathbf{y}_0 because there is no historical data. The use of
 193 recursive posterior predictive distributions in ecology for the integration over the utility is
 194 further discussed in Williams et al. [2018] and Leach et al. [2022].

195 To approximate the expected utility, it is common to use Monte Carlo integration:

$$E[U(\mathbf{d})] \approx \frac{1}{R} \sum_{r=1}^R u(\boldsymbol{\theta}^{(r)}, \tilde{\mathbf{y}}_1^{(r)}, \mathbf{d}). \quad (3)$$

196 As it is necessary to draw R datasets and fit R posteriors to approximate the expected
 197 utility for each proposed design, this step is the bottleneck of Bayesian design. This
 198 problem is exacerbated when using an iterative search algorithm such as exchange
 199 algorithm [Royle, 2002]. And, due to the necessity of fitting multiple posteriors for every
 200 update of the design, approximating the posterior with a faster approach such as that of
 201 the Laplace approximation does not necessary alleviate the core computational issues. This
 202 bottleneck will motivate the design strategies discussed below.

203 **3.3 Design strategies**

204 We now describe the design strategies compared in this work. We use three general
 205 approaches: Design-based, model-based, and exchange-based. Design-based strategies are
 206 also referred to as empirical in Henrys et al. [2024], meaning designs are developed based
 207 only on observed properties of the available data or design region. This includes, for
 208 example, the classic space-filling approach. Model-based designs involve algorithmic

209 optimization of a model-dependent quantity obtained after a statistical model is fit to the
210 available data. Exchange-based designs use an exchange algorithm or one of its variants to
211 optimize the design locations iteratively. Though the designs generated via exchange
212 algorithm inherently use a model to optimize the utility or loss function, we make a
213 distinction between these and the "model-based designs" in order to distinguish between
214 designs that perform an iterative search and model-based methods that optimize via a
215 greedy rank-based selection after a single iteration of model-fitting. All the approaches are
216 summarized in Table 1, and more elaborate descriptions are provided in the Supplementary
217 material. In this work, whenever we refer to a model or the resulting posterior, we use a
218 LGCP and its estimation of the latent intensity of the point process. This model is fit
219 using the INLA and SPDE method [Rue et al., 2009, Lindgren et al., 2011] and the inlabru
220 package [Bachl et al., 2019]. All further details regarding implementation are described in
221 the Supplementary materials.

Table 1: Design strategies considered in the simulation study, separated by design-based, model-based, and exchange-based. Full algorithmic details are provided in the Supplement.

Design-based		
Random + space-filling		Randomly select m sites subject to a minimum inter-site distance of 5 km.
Unsampled regions	re-	Randomly select m sites from locations with the fewest existing observations.
Stratified sampling	sam-	Randomly select m sites, stratified between sampled and unsampled regions.
Model-based		
Maximum predictive variance		Select the m sites with the largest posterior predictive variance under the model fit to $\tilde{\mathbf{y}}_1^{(r)}$, targeting sites of high potential uncertainty [Mondain-Monval et al., 2024].
Weighted predictive variance		Weight predictive variance by the inverse number of samples in each site, balancing the selection of uncertain locations with poorly sampled locations [Mondain-Monval et al., 2024]
Rejection sampling	sam-	Randomly propose sites and accept with probability proportional to scaled predicted intensity, including threshold to control concentration in high-intensity areas [Liu and Vanhatalo, 2020].
Additive KL		For each candidate site, augment \mathbf{y}_0 with site-specific observations from $\tilde{\mathbf{y}}_1^{(r)}$ and rank sites by the KL divergence between updated and baseline posteriors.
LOO KL		For each candidate site, start with the full $\tilde{\mathbf{y}}_1^{(r)}$ dataset and remove site-specific observations and rank sites by the KL divergence between the full-data and leave-one-out posteriors.
Weighted LOO		Apply an intensity-based weight during the LOO optimization to favor lower-intensity regions.
Utility emulation (GP)		Use a Gaussian process emulator of the utility surface to identify sites with high predicted utility [Overstall and Woods, 2017]. We use both an Additive and LOO version of this approach.
Exchange-based		
Exchange algorithm	algo-	Iteratively replace selected sites with nearby candidates whenever utility improves [Royle, 2002].
Approximate Coordinate Exchange (ACE)	Ex-	Construct the design using the Gaussian emulation of the utility surface and improve on this initial configuration with ACE stochastic updates [Overstall and Woods, 2017, Buchhorn et al., 2024b]. We implement both LOO and Additive versions.

222 3.4 Citizen scientist preference

223 We also evaluate the effect of weighting the site selection by citizen science preference. To
224 weight the utility, we use the function $\tilde{U}(\mathbf{d}) = \alpha P(\mathbf{d}) + (1 - \alpha)U(\mathbf{d})$, where α balances the
225 unweighted utility and a preference function. We use two different features as a proxy for
226 preference: the distance of each site from the nearest primary road in Belgium and the
227 distance from the user’s current location. To smooth the preference scores we use a
228 Gaussian decay function to transform raw distances $d(s)$ from a given feature of interest:
229 $P(\mathbf{d}) = \exp\left(-\frac{d(s)^2}{2\gamma^2}\right)$, where γ is the median of the distances. Both utilities and preference
230 scores are standardized before combination. Other useful metrics of accessibility and
231 preference are possible, such as classifying different types of landcover as more preferable to
232 others, the cost of each visit for the user, the personal user collection history, or the desired
233 sampling duration.

234 3.5 Evaluation strategy

235 To evaluate the optimized designs, we again use the virtual ecologist approach to mimic
236 the situation in which additional data is collected in the real world. However, the
237 evaluation differs from the design optimization procedure in that we use the known fixed
238 effect parameters, as well as a new realization of the spatial field, to generate uncollected
239 \mathbf{y}_1 data. Then, using the selected sites of each design \mathbf{d} resulting from the respective design
240 optimization approach, data from \mathbf{y}_1 is collected and the working model fit to the updated
241 dataset, $(\mathbf{y}_0, \mathbf{y}_{1,\mathbf{d}})$. This is repeated for 50 simulations to account for uncertainty in the
242 collection of \mathbf{y}_1 . To compare each design approach we use four metrics: Mean Absolute

243 Error (MAE) of the predicted intensity, predictive variance averaged over all sites, Root
244 Mean Square Error (RMSE) over all fixed effect parameter estimates, and the posterior
245 standard deviation (SD) over all fixed effect parameters. These metrics balance predictive
246 power and model uncertainty, giving a relatively broad assessment of performance of the
247 design methodology. We provide additional information regarding the evaluation process in
248 the Supplemental material.

249 **4 Evaluation of design methodology**

250 **4.1 Comparisons of strategies**

251 We now conduct a comparative analysis of the design methodologies described above. We
252 evaluate the designs in conditions when the underlying data generating process leads to
253 both sparse and abundant \mathbf{y}_0 data, and we compare designs of size $m \in \{5, 10, 20, 30, 50\}$.
254 Figure 1 shows the full comparison over strategies, methods, and survey sizes. Figure 2
255 summarizes the boxplot spread per method, and Figure 3 ranks the performance of each
256 method over all metrics.

Design methodology performance

Comparison over 50 simulation runs

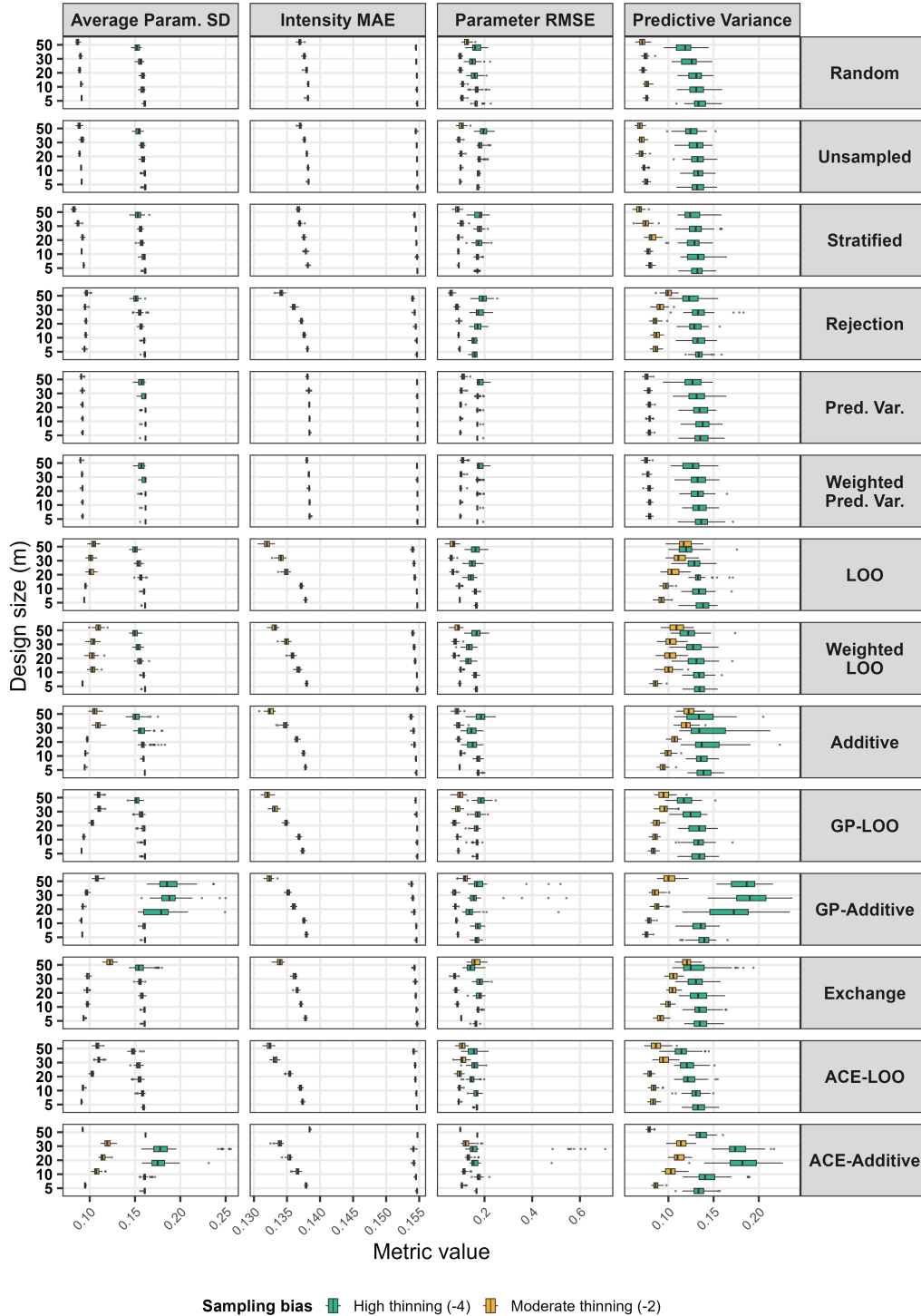


Figure 1: Comparison over methods, metrics, and design survey size. The boxplots show the distribution of performance over 50 simulations of y_1 data collection. Method type is shown on the right y axis; metric type is shown on the top x axis; design size is shown on the left y axis; scores are shown on the bottom x axis. The method types are organized by design-based first, then model-based, then exchange-based. The sparse (high thinning, $\gamma_0 = -4$) and abundant (moderate thinning, $\gamma_0 = -2$) data collection settings are shown within each plot in green and yellow respectively.

257 Figure 1 makes a few design behaviours clear. Across all methods, the performance trends
258 differ quite dramatically between the two levels of thinning, i.e., the richness of available \mathbf{y}_0
259 data. As described in the Supplementary material, to generate \mathbf{y}_0 we thin the data
260 generating process using respective intercepts of $\gamma_0 = -2$ (moderate thinning) and $\gamma_0 = -4$
261 (high thinning), in combination with a road-distance covariate used as a proxy for sampling
262 bias, with the thinning probability parameterized via a logit link function. This
263 construction mirrors collection scenarios with differing degrees of baseline sampling
264 probability for the citizen scientist, and, from the perspective of the design optimization,
265 changes the information available in the starting conditions for the design methods.
266 Consequently, the initial data conditions can lead to very different outcomes when
267 comparing the design-based approaches to the model- and exchange-based approaches.
268 When data is plentiful (thinning= -2), as is often the case with opportunistic data in the
269 big-data paradigm, all approaches improve in metric performance. But the variability
270 between the simulation runs does not necessarily stabilize when more data is available.
271 When considering the spread of the boxplots (the variability of performance across the 50
272 simulations), as more clearly visualized in Figure 2, the design-based methods often exhibit
273 a narrower range of behaviour in the thinning= -2 setting than in the -4 setting. And the
274 design-based methods are much less variable across the simulations when compared to the
275 Additive family (Additive, GP-Additive, and ACE-Additive), especially in the
276 thinning= -4 setting. The LOO family shows less spread than the Additive family in
277 general, but it depends on the method and metric.

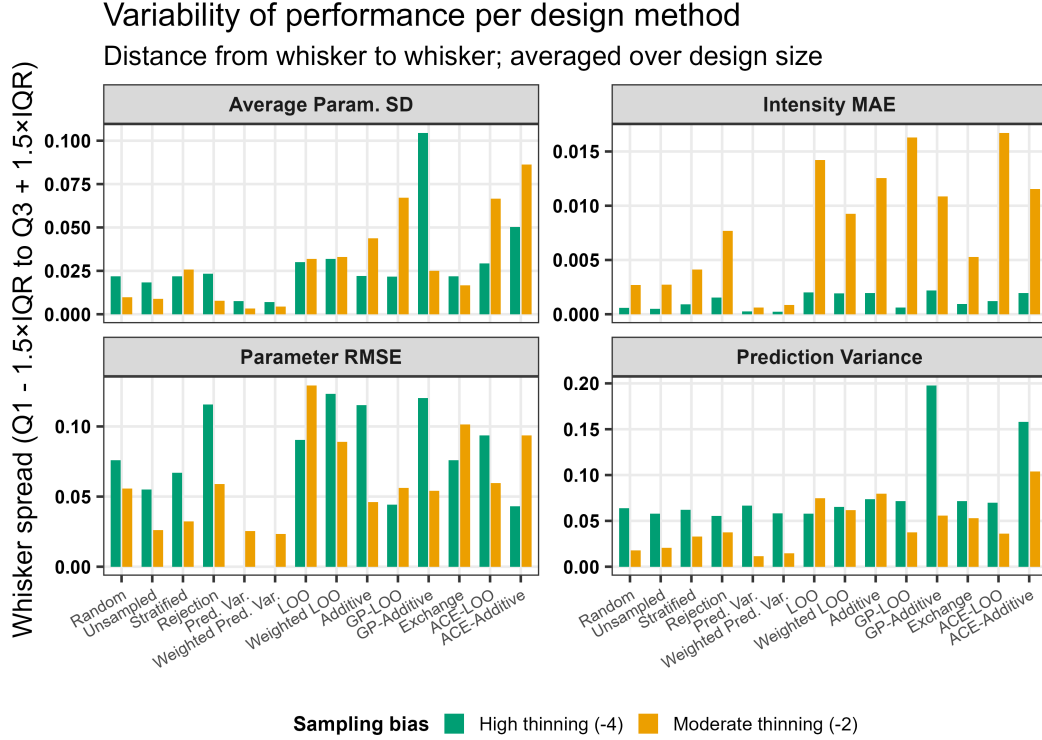


Figure 2: Variability of performance shown in Figure 1 per design method. Bar heights represent the total whisker span, calculated as the distance between $Q_1 - 1.5 \times IQR$ and $Q_3 + 1.5 \times IQR$. The distances are averaged over all design sizes (m) for each metric and method. The thinning conditions are shown for each method.

278 Figure 1 also shows that the Rejection method and the LOO and Additive families improve
 279 in performance of intensity MAE and parameter RMSE as the size of the design increases,
 280 primarily in the abundant thinning= -2 data setting. There are also some improvements
 281 as design size increases in certain thinning= -4 scenarios (e.g., the predictive variance of
 282 GP-LOO and ACE-LOO, and the parameter SD of the entire LOO family). It can also be
 283 seen in Figure 1 that as the design size increases, the median parameter SD and predictive
 284 variance both trend lower in the thinning= -4 setting, but trend higher when
 285 thinning= -2 . This is likely due to conflict between the prior \mathbf{y}_0 data and the imperfectly
 286 simulated $\tilde{\mathbf{y}}_1$ datasets during the design optimization in the thinning= -2 setting. We will
 287 address this further issue in the discussion.

288 The overall performance of each method is further highlighted in Figure 3, where the
289 methods are ranked by the total number of "wins", i.e., the best performance per metric
290 and per simulation of $\tilde{\mathbf{y}}_1$ data collection. In terms of MAE and Parameter RMSE (for
291 thinning= -2), LOO and GP-LOO outperform the other approaches. All three
292 design-based approaches (Random, Unsampled, and Stratified) dominate in terms of
293 parameter SD and predictive variance. Then, in the sparse (thinning= -4) setting, we see
294 a very different ranking. ACE-LOO, Additive, Weighted LOO, GP-Additive are all top
295 performers, each with somewhat different performance profiles. Interestingly, the
296 ACE-LOO performance is composed of parameter SD, RMSE, and predictive variance,
297 suggesting that it may offer the most balanced performance under sparse-data conditions.
298 We see that the Additive family performs well in terms of MAE in the sparse setting, but
299 these methods have very large boxplot spread (Figure 2). However, the strong thinning
300 effect leads to a low-density intensity surface and a compressed error scale, and we
301 therefore don't observe much absolute MAE difference between the methods (see the
302 Intensity MAE column in Figure 1 for thinning= -4).

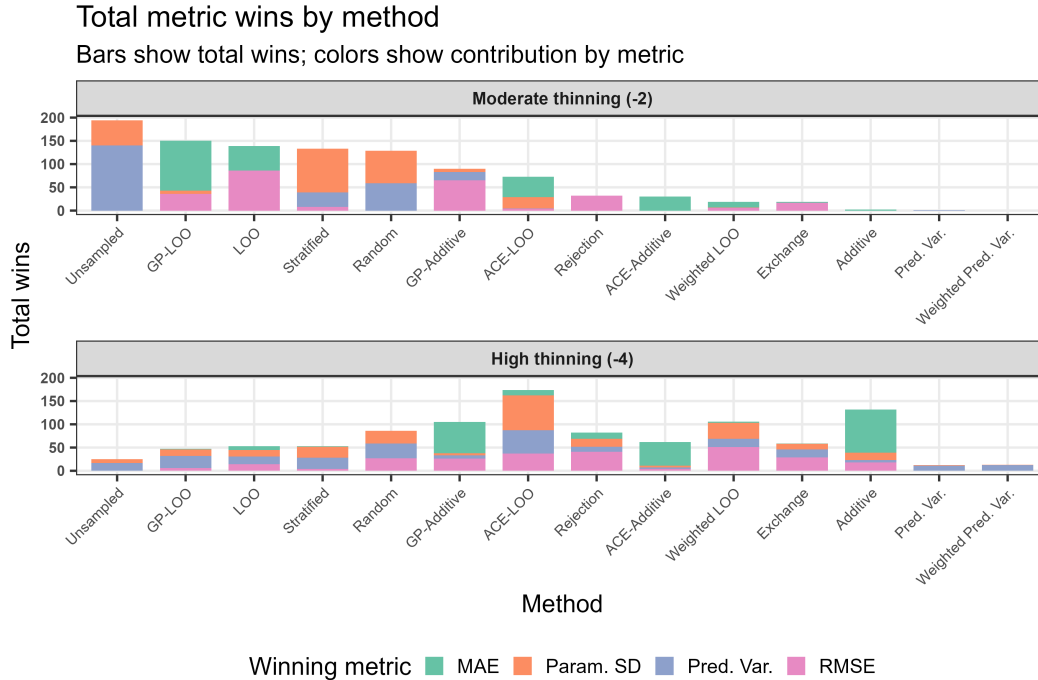


Figure 3: Cumulative performance ranking of design methods, per metric wins. Method order is set by the ranking in the moderate thinning setting. Each bar represents the frequency with which a method achieved the best (lowest) performance for a specific metric per each simulation runs and design size (m). The stacked segments illustrate the contribution of four performance dimensions: Intensity MAE, parameter RMSE, predictive variance, and average parameter SD. Higher total bar heights indicate design methods that consistently outperformed competitors across the evaluation criteria.

303 Finally, in the Supplementary material we include an evaluation of the sensitivity of the
 304 design methods to the number of PPD samples used in the approximation of the expected
 305 utility. The methods tested there showed robustness to the number of PPD samples used,
 306 and between 8-16 samples for most methods was sufficient.

307 4.2 Citizen science preference

308 We next repeat the analysis, now weighting the utilities by citizen science preference. In
 309 these runs we use only the thinning= -2 environment, and we include only a subset of the
 310 methods from the previous analysis: the Additive and LOO families, as well as the

311 predictive and weighted predictive variance approach. We do not include the design-based
312 approaches because they do not use a utility function to select sites.

313 The primary takeaway is that the strength of observer preference or convenience can
314 introduce unexpected bias. We can see this when weighting heavily towards the user's
315 preference (high α), where there is a distinct difference between weighting for distance from
316 current location and weighting for proximity from primary roads. Given Belgium's dense
317 road network, surveyors are never far from main roads at the 10 km resolution we use in
318 our study, and therefore weighting heavily towards road accessibility maintains some
319 space-filling properties of the design. This means that even when heavily weighting towards
320 roads, there is often not a dramatic change in performance of the method. Weighting by
321 location behaves otherwise. As the user's location we use the city of Ghent, an area in
322 which the covariates are homogeneous and the intensity for the simulated process is
323 relatively low, compared to the hot spots in the Southeast of the country. Therefore, when
324 weighting heavily ($\alpha = 0.8$) for location, we observe an increase in MAE, because the
325 model is receiving less information from \mathbf{y}_1 about high intensity areas. Conversely, other
326 metrics exhibit inconsistent, method-dependent responses to intense preference weighting.
327 For example, LOO and Additive both exhibit decreases in parameter SD and predictive
328 variance when $\alpha = 0.8$, relative to the other α levels; and the parameter RMSE doesn't
329 dramatically change, indicating that by including the lower-intensity areas in the design,
330 the model uncertainty has been reduced. In contrast, GP-LOO and GP-additive show
331 increases in parameter RMSE for both distance and roads preference weighting at $\alpha = 0.8$.
332 Furthermore, we don't see large changes in the predictive variance methods because the
333 weighted utilities still strongly favor the border regions of the domain.

Preference design performance

Effect of weighting by roads and location distance over 50 simulations

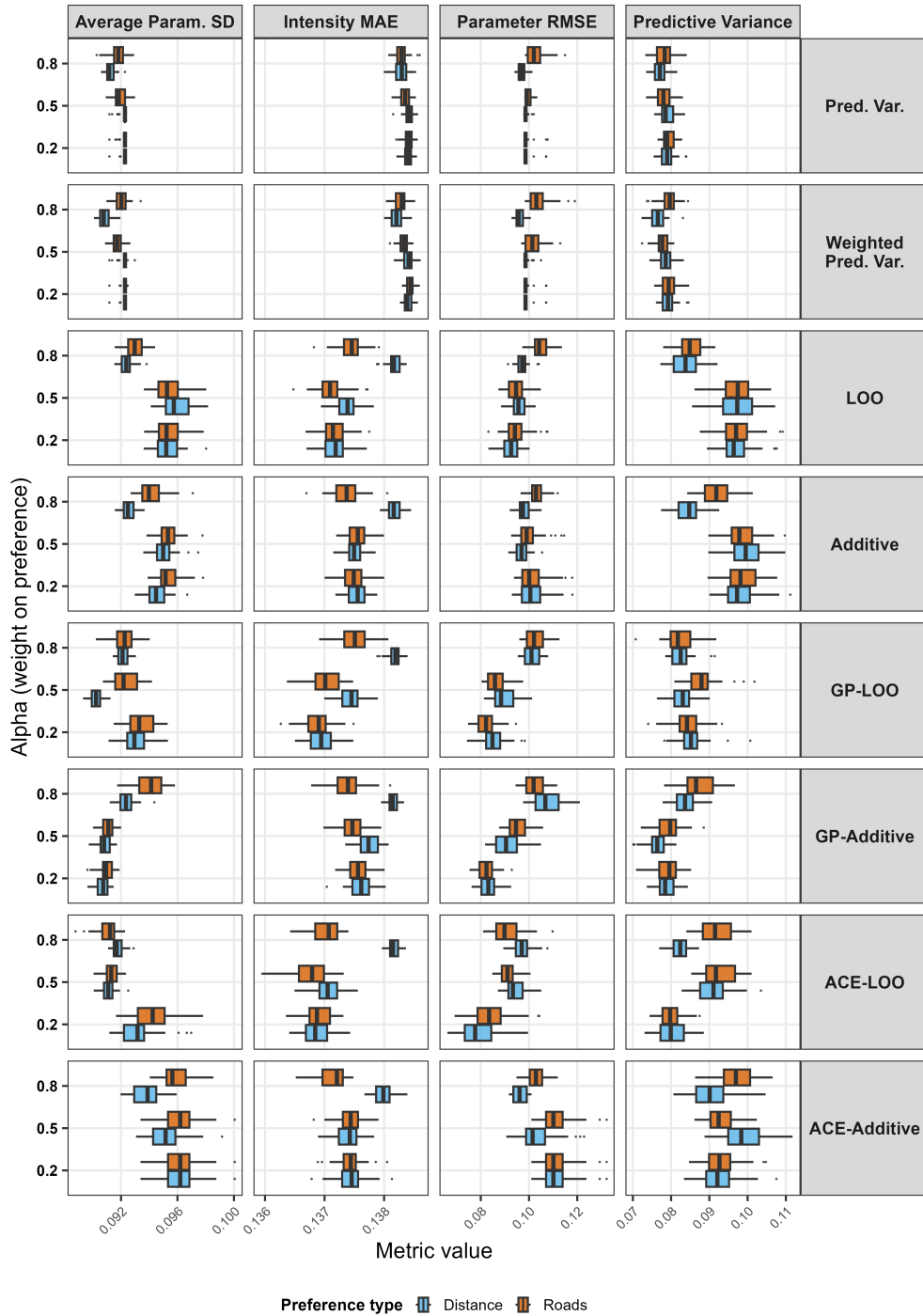


Figure 4: Comparison over methods, metrics, citizen science preference weighting. We only include runs for thinning = -2. The boxplots show the distribution of performance over 50 simulations of \mathbf{y}_1 data collection. Method type is shown on the right y axis; metric type is shown on the top x axis; preference weight (α) is shown on the left y axis; scores are shown on the bottom x axis. The method types are listed by model-based first and then exchange-based. User distance from Ghent (label: Distance) and distance from primary roads (label: Roads) used within the weighting function are shown in blue and orange respectively.

334 As in the analysis without preference weighting, the best method will depend on the
335 optimization goal. Figure 5 shows the ranks of the methods by metric wins. ACE-LOO
336 generally outperforms the other approaches, now in terms of MAE and parameter RMSE.
337 GP-Additive and the Predictive Variance methods outperform the others in terms of
338 parameter SD and predictive variance (which were previously led by the design-based
339 approaches in the unweighted analysis). But because an increase in one metric can lead to
340 a decrease in another, for a real-world application it is important to consider the costs and
341 benefits associated with weighting in the context of the specific design method being used.
342 The core message is then that weighting by preference can have a counterbalancing effect
343 to the purely statistical optimization: Weighting can reduce uncertainty in the model or
344 introduce overconfidence, at the same time driving a directional shift in estimation and
345 prediction accuracy depending on the preferred location and the strength of the weights.

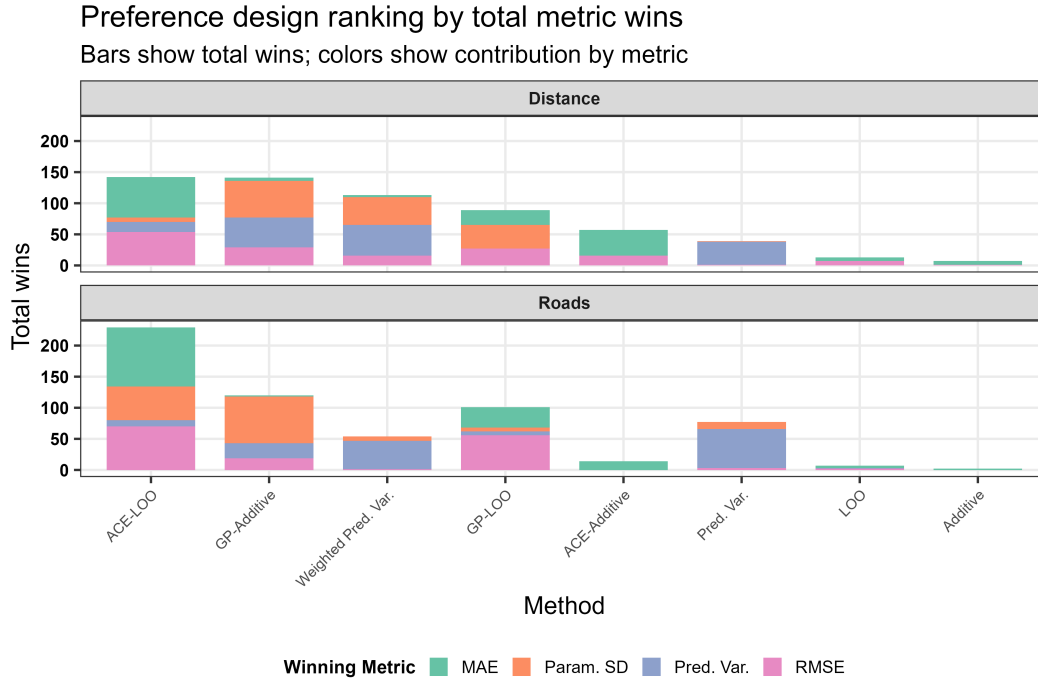


Figure 5: Cumulative performance ranking of design methods, per metric wins. Method order is set by the ranking in the distance from location (Distance) weighting setting. Each bar represents the frequency with which a method achieved the best (lowest) performance per metric across all simulation runs and all α weights. The stacked segments illustrate the contribution of the four performance dimensions: Intensity MAE, parameter RMSE, predictive variance, and average parameter SD. Higher total bar heights indicate design methods that consistently outperformed competitors across the evaluation criteria.

346 5 Discussion

347 We have performed a comprehensive evaluation and methodological benchmark of Bayesian
 348 adaptive design when applied to ecological data collection. We focused specifically on
 349 approaches for opportunistic data collection and the associated technical issues when
 350 developing adaptive designs. While we have focused on presence-only citizen science data,
 351 the framework we use is very flexible and can be applied to many other types of ecological
 352 and environmental data collection. Such data includes presence-absence or
 353 detection-nondetection data, capture-recapture data, animal movement data, phenology

354 data, or continuous data for environmental gradients such as air pollution or water
355 monitoring networks. The issues we consider (design-based vs model-based, sample size,
356 density of prior data, sensitivity of the expected utility approximation, weighting by
357 preference) are general issues anywhere in adaptive sampling of ecological, environmental,
358 and spatial data. We hope that this work encourages further statistical applications for site
359 selection and Bayesian survey design. However, while the evaluation we have performed is
360 extensive, there still remain quite a number of interesting avenues of research as well as a
361 number of limitations to our work. We first offer five high-level key takeaways from this
362 article and then address the limitations.

363 **Takeaway 1: Amount of prior data.** The design approach and the amount of data that
364 has already been collected cannot be disentangled. Though this may appear obvious, it has
365 broad implications for adaptive design. In situations where detection is low or the presence
366 of species is rare, design-based methods can be considered generally safer but less optimal
367 overall. In situations when strong prior knowledge regarding the species can be
368 incorporated into the design, i.e., when previously collected data is abundant, there is more
369 to be gained by using an model-based or algorithmic approach, though at some risk of
370 variability of performance.

371 **Takeaway 2: Metric of optimization.** Project and design decisions are completely
372 conditional on the quantity used for optimization. The more focused the design algorithm
373 is on optimizing a specific quantity, the better it will perform when evaluated in terms of
374 that quantity, at the risk of becoming overly specialized. From this perspective, the exact
375 design method or algorithm is less important than the metric of optimization and
376 evaluation, and the selection of these must match the goals of the project and the species

377 and region of interest.

378 **Takeaway 3: LOO for performance, design-based for stability.** We observed the
379 LOO-family to be consistently well-performing across the evaluation dimensions, and the
380 design-based methods to be generally more consistent across all experimental conditions
381 (prior data availability, design size, and metric). Though we don't endorse the LOO-family
382 as the "best strategy", these methods appear able to balance the tensions between the data
383 and optimization process, even when weighting for citizen science preference. If the project
384 scope is well defined statistically, then the extra computational effort appears worthwhile.
385 Additionally, our results reinforce the idea often found in design literature that bespoke
386 *hybrid* approaches are the most likely to fill specific organizational needs.

387 **Takeaway 4: The computational bottleneck.** The computational bottleneck of
388 Bayesian design should not be intimidating and in fact we observed good stability of the
389 approximation of the expected utility. We find the main complication of model-based
390 design to be the development of the working model, as this will be used to simulate the
391 future \mathbf{y}_1 datasets and must match the data generating process of \mathbf{y}_0 . If the model
392 specification poorly matches the data generating process of \mathbf{y}_0 , it will likely lead to even
393 worse conflicts between the two datasets than we observed during our experiments.

394 **Takeaway 5: Preference weighting:** If taking citizen science preference into account,
395 consider what sort of unintended consequences this might have for the quantity of
396 optimization and how this will interact with the design method of choice. For example, we
397 observed that heavy preference weighting can have a regularizing effect on model
398 parameters: Reducing parameter SD and predictive variance, while simultaneously
399 increasing predictive error and parameter RMSE. When selecting design locations that

400 fulfill stakeholder requirements or qualitative organizational priorities, there will be implicit
401 bias built in to the downstream data collection, even if that bias is not necessarily harmful
402 in terms of model performance.

403 Having provided the high-level messages we seek to communicate with this work, we next
404 discuss a few technical limitations. First, though our working model includes spatial
405 correlation structure, in the model-based design optimization we consider each site
406 independently, in order to avoid the computational complexity of exchange algorithms that
407 compare utilities of m -size configurations. As m and the size of the design space increase
408 (or the resolution of the designs increase), and the complexity of the model increases, there
409 are few available scalable approaches. The work of Gray and Evangelou [2023] considers
410 the site locations jointly, where the authors model the entire sampling design as a random
411 variable with probability proportional to its utility, incorporating experimenter preference
412 into the utility.

413 Moreover, efficient approximations to the expectation of the utility remains the key issue
414 for Bayesian design. Specifically, efficient approaches to the outer loop bottleneck of the
415 expected utility are required, such as via the Sequential Monte Carlo approach of
416 Senarathne et al. [2020], though the application of their method to ecological survey design
417 is not straightforward. Other sequential updating approaches, such as that of Leach et al.
418 [2022], can also improve efficiency by avoiding the full recomputation of the posterior at
419 every additional design point.

420 As mentioned previously, we observed an increase in model uncertainty as the design size
421 increased in the thinning= -2 environment. This may reflect noise in the \mathbf{y}_0 dataset and
422 the increasing contribution from $\mathbf{y}_{1,d}$ to the updated data $(\mathbf{y}_0, \mathbf{y}_{1,d})$ as the design size

423 grows. More generally, simulation of $\tilde{\mathbf{y}}_1$ is a standard step of Bayesian adaptive design
424 when the expected utility is approximated by averaging over posterior predictive draws. As
425 a result, uncertainty in the posterior distribution based on \mathbf{y}_0 propagates into the
426 simulated $\tilde{\mathbf{y}}_1$ values and, consequently, into the utility estimates. Care must therefore be
427 taken to develop a reliable working model, including providing as much relevant prior
428 information as possible, just as in any other application of Bayesian statistical analysis.

429 There are additional ecological limitations of our work. While Mondain-Monval et al.
430 [2024] considers multiple species in their modelling, we use a single simulated species here
431 for the purpose of focusing on the design behaviour itself, unconditional on the species of
432 interest. Tailoring adaptive design approaches to specific species is essential to bring our
433 work closer to ecological reality, as no method offers a one-size-fits-all solution. The
434 effective information content of the prior data depends on species ecology, detectability,
435 spatial structure, and sampling bias inherent to the species. Considering such factors in
436 real-world data collection—in combination with the tensions between data and design
437 methodology that we have explored here—remains a required area of study. Moreover,
438 testing the methods in settings where occupancy probabilities over multiple collection
439 periods have been estimated is essential, in order to be able to compare the design
440 performance to baseline ecological knowledge.

441 Finally, to incorporate Bayesian adaptive design into real-world site selection and sampling
442 projects, the design approaches need to be considered under the full scope of organizational
443 and stakeholder objectives. In this paper we have focused primarily on the statistical
444 aspects of the designs, but we welcome any future collaboration that works towards
445 embedding Bayesian design methods within citizen science projects. We hope to achieve

446 the goal of real-time design described by Callaghan et al. [2021]: "If there are further
447 developments of semi-automated or automated pipelines to interact with citizen scientists
448 in near real time then the collective effort of citizen scientists will be able to reduce
449 redundancies and gaps in the data collected." Integrating the lessons from our work into
450 the Routine-Opportunistic Adaptive Monitoring framework of Pollock et al. [2026] appears
451 to be a clear path forward. Bayesian adaptive design naturally fills the role of their inner
452 loop of opportunistic burst sampling, and potentially can directly communicate with the
453 long-term objectives of the outer monitoring loop. From this perspective, Bayesian
454 adaptive design has the potential to support rapid biodiversity monitoring and
455 environmental assessment under climate and anthropogenic change, while contributing to a
456 more unified ecological framework that balances both statistical ecological priorities.

457 **6 Acknowledgements**

458 This paper received funding from the Flemish Government under the research program:
459 Artificiele Intelligentie (AI) Vlaanderen. Additional thanks are given to the members of the
460 Biostat group in the Faculty of Bioscience Engineering at Ghent University for their helpful
461 conversations related to the statistical topics in this work.

462 **7 Conflict of Interest statement**

463 The authors declare no conflict of interest.

464 **8 Author Contributions**

465 MS and SL conceived the initial ideas for the project. MS implemented the design
466 algorithms and associated code. MS managed the experiments and analyzed the outcomes,
467 and SL provided guidance during this process; MS drafted and wrote the manuscript. SL
468 provided feedback and editing of the manuscript drafts. All authors contributed critically
469 to the drafts and gave final approval for publication.

470 **9 Data Availability**

471 All data used in this work can be found in the anonymous figshare repository at
472 <https://figshare.com/s/61fe2afc676a2a725531>. Following publication of this article in
473 a peer-reviewed journal, the figshare repository will be made public at DOI
474 10.6084/m9.figshare.31915401, along with a GitHub repository for public use of the code
475 and data.

476 **References**

477 Tom August, Richard Fox, David B. Roy, and Michael J. O. Pocock. Data-derived metrics
478 describing the behaviour of field-based citizen scientists provide insights for project
479 design and modelling bias. *Scientific Reports*, 10(1):11009, July 2020. ISSN 2045-2322.
480 doi: 10.1038/s41598-020-67658-3.

481 Fabian E. Bachl, Finn Lindgren, David L. Borchers, and Janine B. Illian. Inlabru: An R

482 package for Bayesian spatial modelling from ecological survey data. *Methods in Ecology*
483 *and Evolution*, 10(6):760–766, 2019. ISSN 2041-210X. doi: 10.1111/2041-210X.13168.

484 Katie Buchhorn, Kerrie Mengersen, Edgar Santos-Fernandez, Erin E Peterson, and
485 James M McGree. Bayesian design with sampling windows for complex spatial processes.
486 *Journal of the Royal Statistical Society Series C: Applied Statistics*, 73(2):378–397,
487 March 2024a. ISSN 0035-9254. doi: 10.1093/jrsssc/qlad099. URL
488 <https://doi.org/10.1093/jrsssc/qlad099>.

489 Katie Buchhorn, Kerrie Mengersen, Edgar Santos-Fernandez, Erin E Peterson, and
490 James M McGree. Bayesian design with sampling windows for complex spatial processes.
491 *Journal of the Royal Statistical Society Series C: Applied Statistics*, 73(2):378–397,
492 March 2024b. ISSN 0035-9254. doi: 10.1093/jrsssc/qlad099.

493 Corey T. Callaghan and Dale E. Gawlik. Efficacy of eBird data as an aid in conservation
494 planning and monitoring. *Journal of Field Ornithology*, 86(4):298–304, 2015. ISSN
495 1557-9263. doi: 10.1111/jfo.12121.

496 Corey T. Callaghan, Alistair G. B. Poore, Richard E. Major, Jodi J. L. Rowley, and
497 William K. Cornwell. Optimizing future biodiversity sampling by citizen scientists.
498 *Proceedings of the Royal Society B: Biological Sciences*, 286(1912):20191487, October
499 2019. doi: 10.1098/rspb.2019.1487.

500 Corey T Callaghan, Alistair G B Poore, Thomas Mesaglio, Angela T Moles, Shinichi
501 Nakagawa, Christopher Roberts, Jodi J L Rowley, Adriana VergÉs, John H Wilshire,
502 and William K Cornwell. Three Frontiers for the Future of Biodiversity Research Using

503 Citizen Science Data. *BioScience*, 71(1):55–63, January 2021. ISSN 0006-3568. doi:
504 10.1093/biosci/biaa131.

505 Corey T. Callaghan, Diana E. Bowler, Shane A. Blowes, Jonathan M. Chase, Mitchell B.
506 Lyons, and Henrique M. Pereira. Quantifying effort needed to estimate species diversity
507 from citizen science data. *Ecosphere*, 13(4):e3966, 2022. ISSN 2150-8925. doi:
508 10.1002/ecs2.3966.

509 Corey T Callaghan, Maureen Thompson, Adam Woods, Alistair G B Poore, Diana E
510 Bowler, Fabrice Samonte, Jodi J L Rowley, Nadiah Roslan, Richard T Kingsford,
511 William K Cornwell, and Richard E Major. Experimental evidence that behavioral
512 nudges in citizen science projects can improve biodiversity data. *BioScience*, 73(4):
513 302–313, April 2023. ISSN 0006-3568. doi: 10.1093/biosci/biad012.

514 Robert M. Dorazio. Accounting for imperfect detection and survey bias in statistical
515 analysis of presence-only data. *Global Ecology and Biogeography*, 23(12):1472–1484, 2014.
516 ISSN 1466-8238. doi: 10.1111/geb.12216.

517 Ian Flint, Chung-Huey Wu, Roozbeh Valavi, Wan-Jyun Chen, and Te-En Lin. Maximising
518 the informativeness of new records in spatial sampling design. *Methods in Ecology and*
519 *Evolution*, 15(1):178–190, 2024. ISSN 2041-210X. doi: 10.1111/2041-210X.14260.

520 Elizabeth J Gray and Evangelos Evangelou. A design utility approach for preferentially
521 sampled spatial data. *Journal of the Royal Statistical Society Series C: Applied*
522 *Statistics*, 72(4):1041–1063, September 2023. ISSN 0035-9254, 1467-9876. doi:
523 10.1093/jrssc/qlad040.

524 Peter A. Henrys, Thomas O. Mondain-Monval, and Susan G. Jarvis. Adaptive sampling in

525 ecology: Key challenges and future opportunities. *Methods in Ecology and Evolution*, 15
526 (9):1483–1496, 2024. ISSN 2041-210X. doi: 10.1111/2041-210X.14393.

527 Mevin B. Hooten, Johnson , Devin S., and Brian M. and Brost. Making Recursive
528 Bayesian Inference Accessible. *The American Statistician*, 75(2):185–194, May 2021.
529 ISSN 0003-1305. doi: 10.1080/00031305.2019.1665584.

530 Alison Johnston, Eleni Matechou, and Emily B. Dennis. Outstanding challenges and future
531 directions for biodiversity monitoring using citizen science data. *Methods in Ecology and*
532 *Evolution*, 14(1):103–116, 2023. ISSN 2041-210X. doi: 10.1111/2041-210X.13834.

533 Vira Koshkina, Yan Wang, Ascelin Gordon, Robert M. Dorazio, Matt White, and Lewi
534 Stone. Integrated species distribution models: Combining presence-background data and
535 site-occupancy data with imperfect detection. *Methods in Ecology and Evolution*, 8(4):
536 420–430, 2017. ISSN 2041-210X. doi: 10.1111/2041-210X.12738.

537 Clinton B. Leach, Perry J. Williams, Joseph M. Eisaguirre, Jamie N. Womble, Michael R.
538 Bower, and Mevin B. Hooten. Recursive Bayesian computation facilitates adaptive
539 optimal design in ecological studies. *Ecology*, 103(2):e03573, 2022. ISSN 1939-9170. doi:
540 10.1002/ecy.3573.

541 Finn Lindgren, Håvard Rue, and Johan Lindström. An explicit link between Gaussian
542 fields and Gaussian Markov random fields: The stochastic partial differential equation
543 approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 73
544 (4):423–498, 2011. ISSN 1467-9868. doi: 10.1111/j.1467-9868.2011.00777.x.

545 Jia Liu and Jarno Vanhatalo. Bayesian model based spatiotemporal survey designs and

546 partially observed log Gaussian Cox process. *Spatial Statistics*, 35:100392, March 2020.
547 ISSN 2211-6753. doi: 10.1016/j.spasta.2019.100392.

548 Quan Long, Marco Scavino, Raúl Tempone, and Suojin Wang. Fast estimation of expected
549 information gains for Bayesian experimental designs based on Laplace approximations.
550 *Computer Methods in Applied Mechanics and Engineering*, 259:24–39, June 2013. ISSN
551 0045-7825. doi: 10.1016/j.cma.2013.02.017. URL
552 <https://www.sciencedirect.com/science/article/pii/S0045782513000492>.

553 Jorge Mateu and Werner G. Muller. *Spatio-Temporal Design: Advances in Efficient Data*
554 *Acquisition*. John Wiley & Sons, Ltd, 2013.

555 Thomas Mondain-Monval, Michael Pocock, Simon Rolph, Tom August, Emma Wright, and
556 Susan Jarvis. Adaptive sampling by citizen scientists improves species distribution model
557 performance: A simulation study. *Methods in Ecology and Evolution*, 15(7):1206–1220,
558 2024. ISSN 2041-210X. doi: 10.1111/2041-210X.14355.

559 Antony M. Overstall and David C. Woods. Bayesian Design of Experiments Using
560 Approximate Coordinate Exchange. *Technometrics*, 59(4):458–470, October 2017. ISSN
561 0040-1706. doi: 10.1080/00401706.2016.1251495.

562 Antony M. Overstall, James M. McGree, and Christopher C. Drovandi. An approach for
563 finding fully Bayesian optimal designs using normal-based approximations to loss
564 functions. *Statistics and Computing*, 28(2):343–358, March 2018. ISSN 1573-1375. doi:
565 10.1007/s11222-017-9734-x. URL <https://doi.org/10.1007/s11222-017-9734-x>.

566 Michael J. O. Pocock, Mark Logie, Nick J. B. Isaac, Richard Fox, and Tom August. The
567 recording behaviour of field-based citizen scientists and its impact on biodiversity trend

568 analysis. *Ecological Indicators*, 151:110276, July 2023. ISSN 1470-160X. doi:
569 10.1016/j.ecolind.2023.110276.

570 Laura J. Pollock, Pedro Henrique Pereira Braga, Christopher R. Florian, Katherine Hébert,
571 Jenna Kline, R. Patrick Lyon, John T. Van Stan, Sara Beery, Michael E. Dillon,
572 Diego Ellis Soto, Brooke Goodman, Niall Hanan, Marta A. Jarzyna, Justin Kitzes, Anke
573 Kügler, Daniel Mosse, Yiluan Song, and Jeff Larkin. Putting the ‘Adaptive’ in Adaptive
574 Monitoring: From Fast Data to Meaningful Ecological Change. *EcoEvoRxiv*, February
575 2026.

576 Henry T. Reich. Optimal sampling design and the accuracy of occupancy models.
577 *Biometrics*, 76(3):1017–1027, 2020. ISSN 1541-0420. doi: 10.1111/biom.13203.

578 J.A Royle. Exchange algorithms for constructing large spatial designs. *Journal of*
579 *Statistical Planning and Inference*, 100(2):121–134, February 2002. ISSN 03783758. doi:
580 10.1016/S0378-3758(01)00127-6.

581 Håvard Rue, Sara Martino, and Nicolas Chopin. Approximate Bayesian inference for latent
582 Gaussian models by using integrated nested Laplace approximations. *Journal of the*
583 *Royal Statistical Society: Series B (Statistical Methodology)*, 71(2):319–392, 2009. ISSN
584 1467-9868. doi: 10.1111/j.1467-9868.2008.00700.x.

585 Elizabeth G. Ryan, Christopher C. Drovandi, James M. McGree, and Anthony N. Pettitt.
586 A Review of Modern Computational Algorithms for Bayesian Optimal Design.
587 *International Statistical Review*, 84(1):128–154, 2016. ISSN 1751-5823. doi:
588 10.1111/insr.12107. URL

589 <https://onlinelibrary.wiley.com/doi/abs/10.1111/insr.12107>. `_eprint:`
590 <https://onlinelibrary.wiley.com/doi/pdf/10.1111/insr.12107>.

591 S. G. J. Senarathne, C. C. Drovandi, and J. M. McGree. A Laplace-based algorithm for
592 Bayesian adaptive design. *Statistics and Computing*, 30(5):1183–1208, September 2020.
593 ISSN 1573-1375. doi: 10.1007/s11222-020-09938-6.

594 S. G. J. Senarathne, Werner G. Müller, and James M. McGree. Bayesian design for
595 minimizing prediction uncertainty in bivariate spatial responses with applications to air
596 quality monitoring. *Biometrical Journal*, 65(4):2100386, 2023. ISSN 1521-4036. doi:
597 10.1002/bimj.202100386. URL
598 <https://onlinelibrary.wiley.com/doi/abs/10.1002/bimj.202100386>. `_eprint:`
599 <https://onlinelibrary.wiley.com/doi/pdf/10.1002/bimj.202100386>.

600 Emily G. Simmonds, Susan G. Jarvis, Peter A. Henrys, Nick J. B. Isaac, and Robert B.
601 O’Hara. Is more data always better? A simulation study of benefits and limitations of
602 integrated distribution models. *Ecography*, 43(10):1413–1422, 2020. ISSN 1600-0587. doi:
603 10.1111/ecog.05146.

604 Daniel Simpson, Håvard Rue, Andrea Riebler, Thiago G. Martins, and Sigrunn H. Sørbye.
605 Penalising Model Component Complexity: A Principled, Practical Approach to
606 Constructing Priors. *Statistical Science*, 32(1):1–28, February 2017. ISSN 0883-4237,
607 2168-8745. doi: 10.1214/16-STS576. URL <https://projecteuclid.org/journals/statistical-science/volume-32/issue-1/Penalising-Model-Component-Complexity--A-Principled-Practical-Approach-to/10.1214/16-STS576.full>. Publisher:
608 Institute of Mathematical Statistics.

611 A. W. L. Pubudu Thilan, Erin Peterson, Patricia Menéndez, Julian Caley, Christopher
612 Drovandi, Camille Mellin, and James McGree. Bayesian design methods for improving
613 the effectiveness of ecosystem monitoring. *Environmental and Ecological Statistics*, 31(4):
614 893–919, December 2024. ISSN 1573-3009. doi: 10.1007/s10651-024-00623-9.

615 Awlp Thilan, P Menéndez, and Jm McGree. Assessing the ability of adaptive designs to
616 capture trends in hard coral cover. *Environmetrics*, 34(6):e2802, 2023. ISSN 1099-095X.
617 doi: 10.1002/env.2802.

618 Maureen M. Thompson, Katie Moon, Adam Woods, Jodi J. L. Rowley, Alistair G. B.
619 Poore, Richard T. Kingsford, and Corey T. Callaghan. Citizen science participant
620 motivations and behaviour: Implications for biodiversity data coverage. *Biological
621 Conservation*, 282:110079, June 2023. ISSN 0006-3207. doi:
622 10.1016/j.biocon.2023.110079.

623 Perry J. Williams, Mevin B. Hooten, Jamie N. Womble, George G. Esslinger, and
624 Michael R. Bower. Monitoring dynamic spatio-temporal ecological processes optimally.
625 *Ecology*, 99(3):524–535, 2018. ISSN 1939-9170. doi: 10.1002/ecy.2120.

626 Elise F Zipkin, Erin R Zylstra, Alexander D Wright, Sarah P Saunders, Andrew O Finley,
627 Michael C Dietze, Malcolm S Itter, and Morgan W Tingley. Addressing data integration
628 challenges to link ecological processes across scales. *Frontiers in Ecology and the
629 Environment*, 19(1):30–38, 2021. ISSN 1540-9309. doi: 10.1002/fee.2290.

630 Damaris Zurell, Uta Berger, Juliano S. Cabral, Florian Jeltsch, Christine N. Meynard,
631 Tamara Münkemüller, Nana Nehrbass, Jörn Pagel, Björn Reineking, Boris Schröder, and

632 Volker Grimm. The virtual ecologist approach: Simulating data and observers. *Oikos*,
633 119(4):622–635, 2010. ISSN 1600-0706. doi: 10.1111/j.1600-0706.2009.18284.x.

634 **Supplementary Material**

635 **S1 Virtual environment and data simulation**

636 **Study region**

637 To evaluate the adaptive design framework, we construct a virtual environment using the
638 geographic extent of Belgium as a case-study. The study domain D was defined using
639 official Belgian administrative boundaries (Lambert 72 projection, EPSG:3137). To ensure
640 a clean simulation window, we processed the national border by removing small slivers or
641 internal exclaves within the Netherlands and Germany, as well as applying a small concave
642 hull adjustment to smooth boundary irregularities. We also scaled the coordinate system
643 from meters to kilometers to improve the numerical stability of the LGCP parameter
644 estimation. The continuous region D was discretized into a regular grid of c cells, with a
645 user-defined resolution (10 km in the experiments shown in the paper). Each cell is referred
646 to as site A_i . This grid serves as the basis for covariate aggregation and the subsequent
647 intensity surface calculation.

648 **Covariate processing**

649 We incorporated four primary environmental drivers to represent the true process
650 underlying species distribution and sampling bias:

- 651 • Vegetation: We used the Enhanced Vegetation Index (EVI) MODIS product

652 (MCD12Q2 Version 6.1).

- 653 • Topography: Elevation data was obtained from the NASADEM Digital Elevation
654 Model.
- 655 • Land cover: Categorical land cover data (CORINE Land Cover 2018) was processed
656 to calculate the fractional cover of specific classes within each grid cell. We use
657 fractional covariates of the classes "2 - Discontinuous urban fabric", "12 -
658 Non-irrigated arable land", and "24 - Coniferous forest". For the purpose of our
659 simulation study these classes serve as representative proxies for the primary
660 environmental and anthropogenic drivers of species intensity within the study region.
- 661 • Accessibility features: To simulate accessibility bias, we calculated the Euclidean
662 distance from the centroid of each grid cell to the nearest primary highway. The
663 highway data was obtained from the Top10Vector (v7.0, 2023) dataset maintained by
664 the National Geographic Institute of Belgium.

665 All covariates were standardized prior to simulation.

666 Point process simulation

667 The true latent species distribution is represented by a log-intensity surface $\log(\lambda(s))$. This
668 data generating model is defined as a linear combination of environmental fixed effects and
669 a latent spatial random field:

$$\log(\lambda(s)) = \beta_0 + \sum_{k=1}^K \beta_k X_k(s) + \omega(s) \quad (\text{S1})$$

670 where β_k are the true regression parameters for the k -th environmental covariate X_k . $\omega(s)$
671 is a Gaussian Random Field (GRF) with Matérn covariance structure simulated via the
672 package `gstat`. The GRF represents latent environmental drivers or inherent properties of
673 the population process that influence species distribution beyond the remote-sensing
674 covariates. The Matérn covariance was configured with smoothness parameter $\nu = 1$,
675 variance $\sigma^2 = 0.5$, and spatial decay parameter $\phi = 5$ km.

676 The resulting intensity surface was then sampled to produce the set of points
677 $\mathbf{y}_0 = \{s_1, s_2, \dots, s_n\}$, representing the observed historical data upon which the adaptive
678 sampling designs are optimized. The generation of \mathbf{y}_0 follows a two-step process. First, a
679 realization of the Log-Gaussian Cox Process (LGCP) is generated based on the latent
680 intensity $\lambda(s)$: $\mathbf{y}_{\text{init}} \sim \text{Poisson Process}(\lambda(s))$ using the `spatstat` package. Next, the
681 probability that a true individual at location s is observed, with probability $p(s_i)$, is
682 constructed via a logit link function using bias-inducing covariates (e.g., proximity to
683 roads): $\text{logit}(p(s_i)) = \gamma_0 + \sum_{j=1}^J \gamma_j B_j(s)$ where $B_j(s)$ represents the j -th bias covariate and
684 γ_j is its associated effect. Each raw point is subjected to a Bernoulli trial with probability
685 $p(s_i)$. Points that fail the trial are discarded (thinned), resulting in the final observed
686 dataset \mathbf{y}_0 . This two-stage data generation procedure decouples the biological process
687 (fixed effects β s) from the sampling process (bias effects γ) in the data generating
688 procedure. In this work we use the distance from primary highways in Belgium as the bias
689 covariate, in addition to an intercept, meaning that as the distance increases, the
690 probability of observation decreases when the parameter associated with the covariate is
691 negative. In our simulations we use two different starting conditions for the bias: $\gamma_0 = -2$
692 and $\gamma_0 = -4$. The -4 setting generates a much higher rate of thinning ("high" thinning in

693 the main text), i.e., an overall lower baseline probability of making a positive observation.
 694 This could be attributed to a number of reasons: For example, a difficult to detect species,
 695 a rare species, or a species of less interest to volunteer observers. We include a map of this
 696 thinning probability in Supplemental figure S1, and the realizations of the data generating
 697 process in Supplemental figure S2.

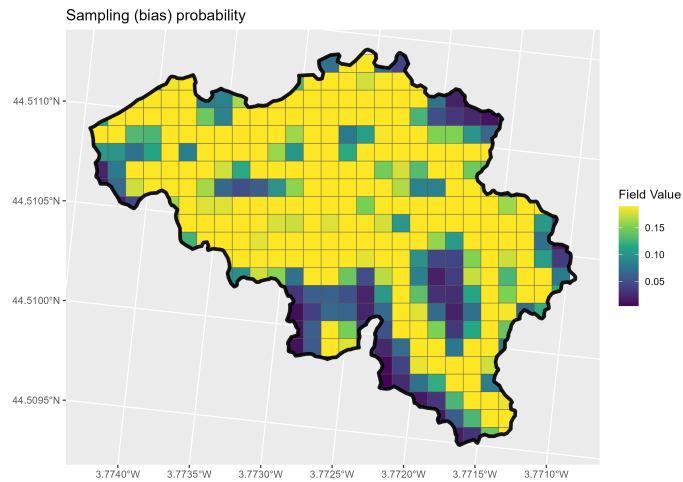


Figure S1: Probability of making a positive observation of the simulated species within Belgium, $\gamma_0 = -2$, based on the parameterized thinning process using the logit link. The overall probability is relatively low which reflects the inherent sparsity and incomplete coverage that is characteristic of opportunistic datasets.

698 S2 Modelling

699 To estimate the species distribution and recover the environmental parameters and latent
 700 spatial field, we apply the LGCP model of Supplemental equation S1. However, in the
 701 working model the probability of observation is parameterized via a log-linear relationship,
 702 intentionally introducing a source of misspecification. We fit the model using the INLA
 703 (Integrated Nested Laplace Approximation) and Stochastic Partial Differential Equation

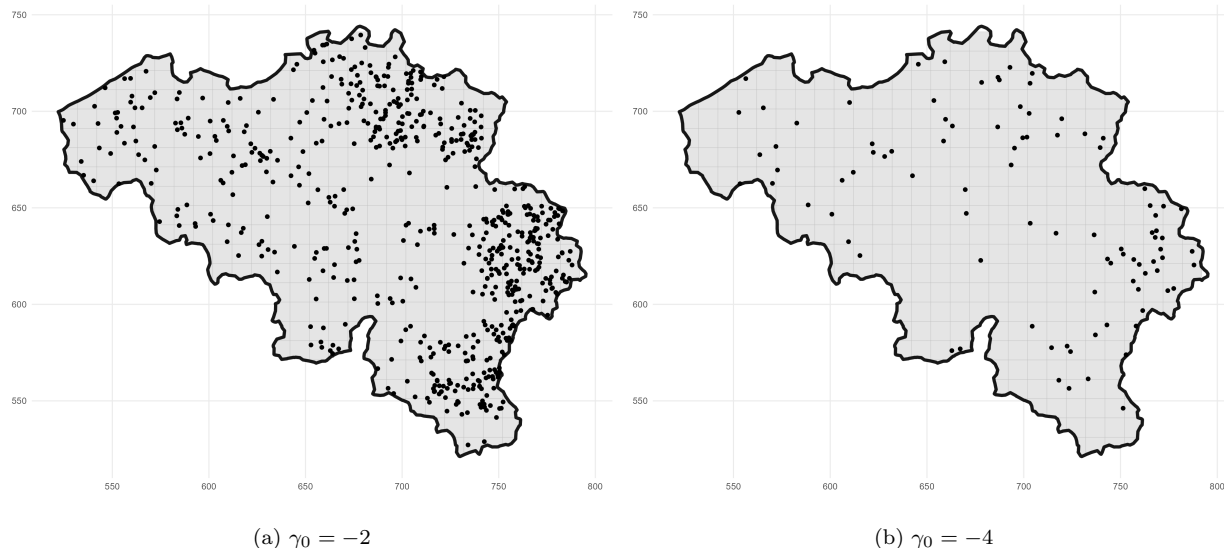


Figure S2: Simulated \mathbf{y}_0 PO datasets in Belgium resulting from data generating process, after thinning has been performed using the parameterized thinning link. The black points represent locations at which an opportunistic observation has been collected. (a) shows the data from which fewer points have been thinned ($\gamma_0 = -2$ or "moderate" thinning), and (b) shows the data where the thinning effect is much stronger ($\gamma_0 = -4$ or "high" thinning).

704 (SPDE) approach [Rue et al., 2009, Lindgren et al., 2011]. For the Matérn covariance, we
 705 use Penalized Complexity (PC) priors [Simpson et al., 2017]: Spatial range:
 706 $P(\rho < 5 \text{ km}) = 0.01$. Marginal standard deviation: $P(\sigma > 0.5) = 0.01$. The data
 707 generating model uses the gstat range ϕ , which is related to the SPDE range by
 708 $\rho = \phi \cdot \sqrt{8\nu}$ [Lindgren et al., 2011]. The choice of range, $\phi = 5$, implies a source of
 709 misspecification between the data generating and working models. Because we work with
 710 aggregated 10 km resolution data, the recovery of fine-scale variation is limited, a common
 711 issue in ecological modelling. We induce this misspecification to mimic real-world data
 712 collection and modelling scenarios.

713 Inference was performed using the Laplace approximation strategy of INLA. We use the
 714 `inlabru` package for model fitting. Because the adaptive design requires fitting models
 715 iteratively or in parallel across many posterior predictive datasets, we optimized the

716 computation in the following ways: Integration strategy: We used the Empirical Bayes
717 integration strategy to reduce computational overhead. Priors: The intercept was given a
718 Gaussian prior with a mean of -3 and a precision of 1, reflecting our prior knowledge of the
719 expected point density across the study area. For iterative model updates, we used the
720 control.mode option in inlabru to restart the optimization from the previous model’s
721 configuration.

722 **S3 Utility functions**

723 Within the design optimization, there are a number of interesting utility and loss functions
724 that can be considered. In this paper we use two general types: Kullback-Leibler (KL)
725 divergence and predictive variance, though we compare weighted variants as well, described
726 in the design algorithm section in the main text and in the section below. A more
727 elaborate discussion of utility functions and their properties can be found in Ryan et al.
728 [2016]. Interested readers may consider work such as Overstall and Woods [2017], Overstall
729 et al. [2018], or Liu and Vanhatalo [2020] to see additional examples of the utilities used
730 within Bayesian design.

731 For a design d , the KL utility function can be written as follows, for example when
732 comparing posteriors resulting from \mathbf{y}_1 and \mathbf{y}_0 :

$$u_{\text{KL}}(\mathbf{y}_1, \mathbf{d}) = D_{KL}(p(\boldsymbol{\theta}|\mathbf{y}_1, \mathbf{y}_0, \mathbf{d})|| p(\boldsymbol{\theta}|\mathbf{y}_0)) = \int_{\boldsymbol{\theta}} p(\boldsymbol{\theta} | \mathbf{y}_1, \mathbf{y}_0, \mathbf{d}) \log \left\{ \frac{p(\boldsymbol{\theta} | \mathbf{y}_1, \mathbf{y}_0, \mathbf{d})}{p(\boldsymbol{\theta} | \mathbf{y}_0)} \right\} d\boldsymbol{\theta}. \quad (\text{S2})$$

733 where, in our case, $\boldsymbol{\theta}$ denotes the full set of unknown model quantities including fixed

734 regression effects and latent spatial effects. \mathbf{y}_0 may also be conditional on a previous design
735 \mathbf{d}_0 but we do not consider such prior collection structure here. When estimating the
736 expected utility with MC integration, we must first approximate the posterior of the model
737 and then the expectation of the utility. If using MCMC to estimate the posterior, this
738 means we rely on a computationally expensive double approximation. To avoid the inner
739 expensive approximation, in this work we use the Laplace approximation as implemented
740 via the INLA software. The advantage of using the Laplace approximation, in addition to
741 speed, is that when two multivariate normal distributions are compared, KL divergence can
742 be expressed in the analytic form:

$$u_{\text{KL}}(\mathbf{y}_1, \mathbf{d}) = \frac{1}{2} \left[-\ln \left(\frac{|\Sigma_{y_1}|}{|\Sigma_{y_0}|} \right) + \text{tr}((\Sigma_{y_0})^{-1} \Sigma_{y_1}) + (\theta_{y_0} - \theta_{y_1})^T (\Sigma_{y_0}^{-1}) (\theta_{y_0} - \theta_{y_1}) - k \right]. \quad (\text{S3})$$

743 This is similar or the same as the design formulation as used in Long et al. [2013], Liu and
744 Vanhatalo [2020], Thilan et al. [2024], Senarathne et al. [2023], Buchhorn et al. [2024a] and
745 Overstall et al. [2018], all of whom use the normal distribution and normal form of the KL
746 divergence to approximate the expected loss.

747 The expectation of the utility over the posterior predictive distribution will be

$$\mathbb{E}[U_{\text{KL}}(\mathbf{d}) \mid \mathbf{y}_0] = \int_{\mathbf{y}_1} u_{\text{KL}}(\mathbf{y}_1, \mathbf{d}) p(\mathbf{y}_1 \mid \mathbf{y}_0, \mathbf{d}) d\mathbf{y}_1. \quad (\text{S4})$$

748 We also consider predictive variance in our design optimization. For each candidate site,
749 we compute the expected posterior variance of the intensity, conditional on \mathbf{y}_1 and \mathbf{y}_0 and

750 report the average over all sites:

$$E[\text{Var}(\tilde{\lambda}(s_i)|\mathbf{y}_1, \mathbf{y}_0)|\mathbf{y}_0] = \int_{\mathbf{y}_1} \text{Var}(\tilde{\lambda}(s_i)|\mathbf{y}_1, \mathbf{y}_0)p(\mathbf{y}_1|\mathbf{y}_0)d\mathbf{y}_1. \quad (\text{S5})$$

751 **S4 Design methods**

752 Here we provide more detailed descriptions of the design approaches we use in the paper.

753 **Design-based**

- 754 • Random + space-filling design: Randomly sample m sites. Then for each site, sample
755 j new locations and check if any of these locations satisfy a minimum distance
756 criteria so that no points are within 5 kilometers of the new site.
- 757 • Unsampled regions: Randomly sample m sites from locations with the fewest
758 observations taken.
- 759 • Stratified sampling: Random sample m locations, stratified between sampled and
760 unsampled regions.

761 **Model-based**

- 762 • Maximum predictive variance: Based on Supplemental equation S5, select the m sites
763 that lead to the largest predictive variance, in order to target locations of high
764 uncertainty for future sampling. This is a Bayesian version of the the "uncertainty
765 only" method of Mondain-Monval et al. [2024]

- 766 • Weighted predictive variance: Weight the predictive variances by the inverse of the
767 number of samples in the site. This prioritizes sites where there are few samples and
768 the predictive variance is relatively high. This approach is also used by
769 Mondain-Monval et al. [2024].

- 770 • Rejection sampling design: Sample sites at random and accept these candidates into
771 the design with probability p . The design is grown until m sites have been added. We
772 use a model-based quantity, $p = \min(p_{\max}, \frac{\tilde{\lambda}}{\max(\tilde{\lambda})})$. p_{\max} allows the tuning of the
773 design away from solely high-intensity areas. This is based on Liu and Vanhatalo
774 [2020].

- 775 • Additive design: This method evaluates the marginal information gain of each
776 candidate site A_i , $i \in 1, \dots, n$ by iteratively augmenting \mathbf{y}_0 with the observations
777 $s \in A_i$ from the predictive dataset $\mathbf{y}_1^{(r)}$. The model is refit using this augmented
778 $(\mathbf{y}_0, \mathbf{y}_{1,i})$ dataset, and the KL utility is computed between the original posterior and
779 the updated posterior. The m sites that result in the highest KL divergence are
780 selected. This is inspired by Flint et al. [2024] but adapted for Bayesian design.

- 781 • Leave-one-out (LOO): A removal-style method based on the predictive data \mathbf{y}_1 . For
782 each site A_i , $i \in 1, \dots, n$, a LOO dataset $\mathbf{y}_{1-i}^{(r)}$ is generated from the predictive dataset
783 $\mathbf{y}_1^{(r)}$, where the observations $s \in A_i$ are removed. The KL divergence between the
784 posterior of the full predictive dataset and the LOO dataset is computed and the m
785 sites that result in the highest utility are selected. This can be considered as a sort of
786 Bayesian analogue to the statistical leverage approach of Callaghan et al. [2019].

- 787 • Weighted LOO: As the above method disregards the importance of unsampled sites,
788 this approach weights the LOO utility towards lower intensity areas. This weight is

789 based on the inverse of the predicted intensity with tuning parameter β :
790 $w(s) = (1 + \tilde{\lambda}(s))^{-\beta}$ for location s . The product of the KL divergence and $w(s)$ is
791 taken on a site-by-site basis. As β increases, the contribution of high-intensity sites is
792 increasingly penalized.

- 793 • Utility emulation using Gaussian Process (GP): Based on the first phase of the ACE
794 algorithm of Overstall and Woods [2017] and [Buchhorn et al., 2024a], this Bayesian
795 optimization approach selects the m best sites in terms of utility predicted by a GP
796 emulator fit to the utility surface. It requires first computing the full utility for a
797 subset of sites and then fitting the GP to these utilities. The utility itself is flexible,
798 and here we use both the additive and LOO versions.

799 Exchange-based

- 800 • Exchange algorithm [Royle, 2002]. For a given design \mathbf{d} , exchange each site A_i in the
801 design with each candidate site A'_i in its local neighborhood. Update the design with
802 the candidate if the exchange leads to improved utility. Converge if after a full
803 iteration, no changes are made to the design. Within the exchange we use the
804 additive version of the KL utility to compare candidate designs.
- 805 • ACE (Stochastic exchange): We use a stochastic exchange algorithm similar to
806 Overstall and Woods [2017] and [Buchhorn et al., 2024a], where the utility surface is
807 estimated using a Gaussian process and the design is then created by selecting the
808 sites that best maximize the estimated utility. From this design the points are then
809 exchanged and accepted with probability p as in Phase II of the Overstall and Woods
810 [2017] algorithm. These designs are selected via the KL utility, using both the LOO

811 and additive strategies mentioned above.

812 **S5 Design evaluation strategy**

813 To evaluate each candidate design, we again use a simulation-based framework, in order to
814 compare the designs against a known, ground truth intensity surface and known fixed effect
815 parameters. Mirroring the data generating procedure described above, we first simulate a
816 virtual ecological environment and then sample data from this environment according to
817 the candidate designs. All evaluations were conducted using the same LGCP working
818 model and the designs were compared under 50 repeated realizations of future sampling
819 data.

820 We generate a known spatially varying log-intensity surface across the study domain, using
821 the same combination of spatial covariates and a new realization of the spatial random field
822 described above. We then simulated new realizations of hypothetical species observations
823 based on that underlying intensity. These data were "collected" based on the candidate
824 design, and the newly sampled points were appended to the baseline dataset \mathbf{y}_0 to create
825 an updated dataset $(\mathbf{y}_0, \mathbf{y}_{1,d})$.

826 We then fit the LGCP to this updated dataset. To measure the performance of each design
827 strategy, we extracted the posterior estimates from the fitted LGCP models and calculated
828 several performance metrics against the known simulated truth. The metrics are described
829 below:

- Parameter recovery: We evaluated how well the designs recovered the true underlying

relationships between the intensity and the spatial covariates. We calculated the Root Mean Squared Error (RMSE):

$$\text{RMSE}_\beta = \sqrt{\frac{1}{K} \sum_{k=1}^K (\hat{\beta}_k - \beta_k^{\text{true}})^2},$$

830 where $\hat{\beta}_k$ is the posterior mean.

- 831 • Parameter precision: The posterior variance, $\text{Var}(\hat{\beta}_k)$, was extracted to assess
- 832 estimation precision and reported as the average variance over all parameters.
- 833 • Predictive error: We assessed the model’s ability to map the true species distribution
- 834 by predicting the expected intensity λ across the spatial grid. The predictions were
- 835 evaluated using the Mean Absolute Error (MAE):

$$\text{MAE}_\lambda = \frac{1}{C} \sum_{i=1}^C |\hat{\lambda}_i - \lambda_i^{\text{true}}|,$$

836 where C is the total number of grid cells, $\hat{\lambda}_i = \exp(\hat{\eta}_i + \sigma_i^2/2)$ is the expected

837 intensity for a log-normal variable with predictive variance σ_i^2 , λ_i^{true} is the ground

838 truth intensity at cell i , and $\hat{\eta}_i$ is the log-linear predictor of the intensity.

- Predictive variance: To evaluate overall predictive confidence, we extracted the predictive variance $\hat{\sigma}_i^2$, averaged across all grid cells:

$$\bar{\sigma}_{\text{pred}}^2 = \frac{1}{C} \sum_{i=1}^C \hat{\sigma}_i^2.$$

839 **S6 Sensitivity analysis of Monte Carlo integration**

840 In this section we show the sensitivity of the expected utility to the Monte Carlo
841 approximation. The Monte Carlo approximation of the expected utility is

$$E[U(\mathbf{d})|\mathbf{y}_0] \approx \frac{1}{R} \sum_{r=1}^R u(\boldsymbol{\theta}^{(r)}, \tilde{\mathbf{y}}^{(r)}, \mathbf{d}). \quad (\text{S6})$$

842 We compare design results for increasing number of samples (increasing R), i.e., draws $\tilde{\mathbf{y}}_1^{(r)}$
843 from the Posterior Predictive Distribution and respective posterior fits $\boldsymbol{\theta}^{(r)}$. We focus on
844 only the LOO family, including also predictive and weighted predictive variance, and the
845 exchange algorithm. From Supplemental figure S3, we can see that the performance of the
846 designs begins to converge as R increases, particularly for the LOO family of methods. The
847 predictive and weighted predictive variance methods remains constant across all values of
848 R , as these designs target the boundaries of the design region and are not sensitive to the
849 realizations of \mathbf{y}_1 . The Exchange algorithm shows higher sensitivity in these runs

850 For the LOO family, at the lowest number of integration samples ($R = 1$ to $R = 4$) we
851 observe higher variability across all metrics the Average Parameter SD. This suggests that
852 using a single PPD draw provides an unstable approximation of the expected utility,
853 leading the algorithm to select sub-optimal sites. As R increases, the most dramatic
854 improvements in performance occur from $R = 4$ to $R = 16$. Beyond $R = 24$ the
855 performance begins to plateau, and in general, using approximately 8 to 24 draws appears
856 sufficient to capture the utility signal required for site selection.

Sensitivity of design performance
Effect of PPD integration over 50 simulations

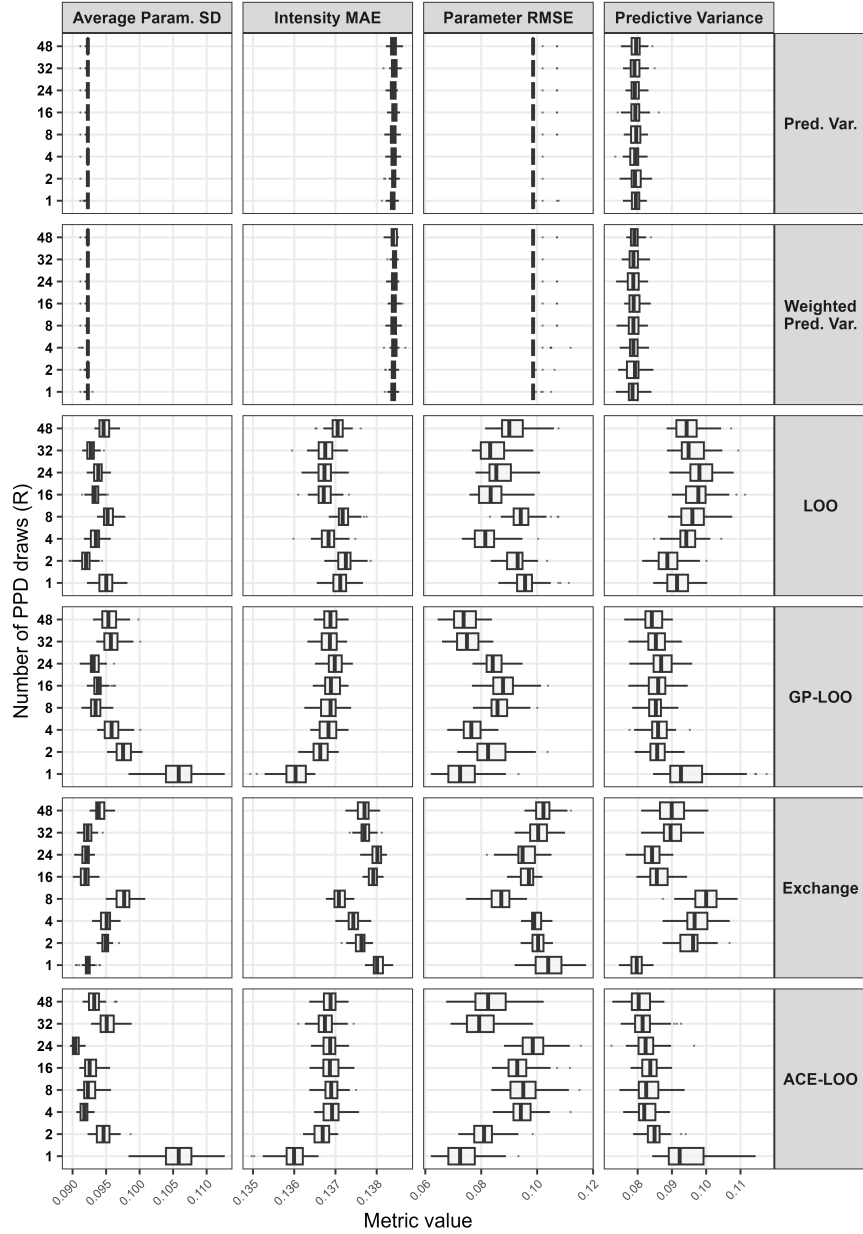


Figure S3: Sensitivity of design performance to the number of Monte Carlo draws (R) from the Posterior Predictive Distribution of \mathbf{y}_0 . The boxplots show the distribution design performance under each metric: Average Parameter SD, Intensity MAE, Parameter RMSE, and Predictive Variance, across 50 simulation runs under the moderate thinning condition ($\gamma_0 = -2$). Rows show the design methods. The y-axis on the left denotes the number of PPD datasets (R) used to approximate the expected utility $U(\mathbf{d})$. Results indicate convergence in performance as R increases, with model- and exchange-based methods exhibiting higher volatility at fewer PPD draws, and the predictive variance methods showing stability across R .