

Combining scientific survey and commercial catch data to map fish distribution

Baptiste Alglave^{1,2*}, Etienne Rivot², Marie-Pierre Etienne³, Mathieu Woillez⁴, James T. Thorson⁵, & Youen Vernard¹

¹ DECOD (Ecosystem Dynamics and Sustainability), IFREMER, Institut Agro, INRAE, Nantes, France

² DECOD (Ecosystem Dynamics and Sustainability), Institut Agro, IFREMER, INRAE, Rennes, France

³ Mathematical Research Institute of Rennes IRMAR, Rennes University, Rennes, France

⁴ DECOD (Ecosystem Dynamics and Sustainability), IFREMER, Institut Agro, INRAE, Brest, France

⁵ Habitat and Ecological Processes Research Program, Alaska Fisheries Science Center, National Marine Fisheries Service, NOAA, Seattle, WA, USA

*corresponding author: email: baptiste.alglave@agrocampus-ouest.fr; present address: DECOD (Ecosystem Dynamics and Sustainability), Institut Agro, IFREMER, INRAE, Rennes, France

Abstract

Developing Species Distribution Models (SDM) for marine exploited species is a major challenge in fisheries ecology. Classical modelling approaches typically rely on fish research surveys data. They benefit from a standardized sampling design and a controlled catchability, but they usually occur once or twice a year and may sample a relatively small number of spatial locations. Spatial monitoring of commercial data (based on logbooks crossed with Vessel Monitoring Systems) can provide an additional extensive data source to inform fish spatial distribution. We propose a spatial hierarchical framework integrating both data sources while accounting for preferential sampling (PS) of commercial data. From simulations, we demonstrate PS should be accounted for in estimation when PS is actually strong. When commercial data far exceed scientific data, the later bring little information to spatial predictions in the areas sampled by commercial data, but bring information in areas with low fishing intensity and provide a validation dataset to assess the integrated model consistency. We applied the framework to three demersal species (hake, sole, and squids) in the Bay of Biscay that emphasize contrasted PS intensity and we demonstrate that the framework can account for several fleet with varying catchabilities and PS behaviors.

Keywords: Species distribution model, integrated modelling, Hierarchical model, VMS and logbook data, survey data, Template Model Builder (TMB)

33 **1 INTRODUCTION**

34 Developing species distribution models (SDM) is critical in marine and fisheries ecology
35 for assessing the relationship between species and their habitat (Guisan and
36 Zimmermann, 2000), identifying essential habitats (Paradinas *et al.*, 2015), forecasting
37 population and ecosystems response to environmental changes (Cheung *et al.*, 2009).
38 The development of statistical models to predict fishery resources distribution has
39 received considerable attention (Planque *et al.*, 2011; Thorson *et al.*, 2015a, 2015b;
40 Martínez-Minaya *et al.*, 2018; Moriarty *et al.*, 2020). Recent developments have
41 generalized SDM to analyze biological data representing condition, stomach contents,
42 size structure, and other demography and population dynamics features (Thorson, 2015;
43 Grüss *et al.*, 2020). Ongoing research also seek to integrate individual movement, growth,
44 species interactions into SDM (Kristensen *et al.*, 2014; Thorson *et al.*, 2017a, 2019),
45 although these approaches are “data hungry” and therefore require integrating different
46 sources of data within a single model.

47 Scientific survey and commercial catch data consist in two potentially complementary data
48 sources to estimate harvested fish spatial distribution (Pennino *et al.*, 2016). Scientific
49 surveys are key data sources in fisheries ecology. They most often benefit from a
50 standardized sampling plan and a constant catchability (Hilborn and Walters, 1992; Ocean
51 Studies Board and National Research Council, 2000; ICES, 2005; Nielsen, 2015). They
52 are generally designed to cover the full geographical extent of specific populations
53 including areas of low or null abundance, and are thus adapted to develop unbiased
54 abundance indices and spatial predictions of species distribution (Rivoirard *et al.*, 2008;
55 ICES, 2012). In addition, they often seek to minimize selectivity in order to sample as
56 many species, size groups and life stages as possible. However, the related expansive

57 charges generally comes at the cost of a relatively low sampling density in space and/or
58 time. For instance, trawl survey can sample a limited number of spatial locations, and
59 most often occur once or twice a year. Thus, they may provide poor information regarding
60 intra-annual variability (Pennino *et al.*, 2016; Rufener *et al.*, 2021) and imprecise estimates
61 of species abundance and distribution (ICES, 2005).

62 Commercial catch declarations (logbooks) data constitute a complementary data source
63 that may benefit of a higher sampling effort than scientific survey. In Europe, catch
64 declarations must be reported in logbooks data for all fishing vessels; besides, geolocation
65 through Vessel Monitoring System (VMS) is mandatory for all fishing boat above 12m long
66 (Hintzen, 2021). Hence, logbook data combined with VMS data can provide high
67 resolution maps of Catch Per Unit Effort (CPUE - Gerritsen and Lordan, 2010; Murray *et*
68 *al.*, 2013) with a relatively dense spatio-temporal sampling within the range of the
69 commercial fleets. However, inferring SDM with commercial data can be challenging as
70 they generally arise from a preferential sampling (PS) behavior, i.e. a sampling that
71 directly or indirectly depends upon the biomass of the target species. Indeed, fishermen
72 tend to target areas with high biomass or may also favor fishing zones based on other
73 criteria (like bottom substrate or distance to the coast for instance - Hintzen *et al.*, 2021)
74 that are indirectly related to the target species abundance. When not properly considered
75 in statistical models, PS associated with commercial data may lead to biased estimates
76 of fish distribution and biomass (Trenkel *et al.*, 2013; Pennino *et al.*, 2019). In particular,
77 when the biomass is spatially heterogeneous, ignoring PS may overestimate the spatial
78 predictions and the overall biomass estimates.

79 Recent research has tackled this challenge and proposed methods to account for PS in
80 statistical inferences. Model based PS was first introduced in a statistical context by Diggle

81 *et al.* (2010). The authors extended a standard geostatistical approach within a
82 hierarchical framework where the variable of interest is modeled as a latent field, with both
83 direct observations and the local intensity of the sampling effort that depend on the latent
84 field. More specifically, the sampling process is modelled as an inhomogeneous Poisson
85 point process whose intensity directly depends on the latent field values. This approach
86 was extended by Pati *et al.* (2011) who introduced covariates and random effects in the
87 model. Conn *et al.* (2017) followed the same ideas and developed a more generic model
88 for ecological applications, which they applied to aerial seal count data. Pennino *et al.*
89 (2019) applied similar ideas to infer the distribution of shrimps from onboard fishery data.
90 Provided PS is accounted for, integrated models (IM) appear as an attractive tool to
91 combine fishery-independent and fishery-dependent data to infer harvested fish spatial
92 distribution. IM have received considerable attention in the ecological literature (Schaub
93 and Abadi, 2011; Parent and Rivot, 2012; Gimenez *et al.*, 2014). By sharing the
94 information between different data types, IM may provide more accurate estimates and
95 predictions compared with separate analysis of different data types. Recently, Rufener *et*
96 *al.* (2021) demonstrated the potential of IM to integrate scientific data and onboard
97 observer count data to improve SDM of fishery resources. However, although onboard
98 observer data provide useful complementary information to scientific survey, they
99 generally only represent a small proportion of all sea trips (1% in average for the French
100 observer programs - Cornou *et al.*, 2021). By contrast, the combination of commercial
101 catch declarations in logbooks with VMS data provides a more extensive data source to
102 map fish spatial distribution. Furthermore, the potential of embedding PS within a
103 hierarchical SDM to integrate catch declaration data and scientific survey is still an open
104 challenge and new methodology are required to handle PS behaviors of commercial fleets

105 while accounting for all the complexity related to fishing locational choice (Salas and
106 Gaertner, 2004; Haynie *et al.*, 2009; Girardin *et al.*, 2017).

107 In this paper, we develop an IM model to infer fish spatial distribution by combining both
108 scientific and commercial catch declaration data while taking into account the PS induced
109 by fishing targeting behavior.

110 To assess the challenges, the benefits but also the limits of the approach, we evaluate
111 the performance of our IM based on simulated data. Simulations are primarily designed
112 to assess the respective contribution of each data source to inference for different model
113 configurations. We first evaluate how the balance between the commercial and scientific
114 sample sizes affect the model outputs. Because the commercial data may often only
115 partially cover the distribution area of a targeted species, we assess how this issue may
116 affect the quality of estimation and how scientific data may contribute to reduce the effect
117 of this gap in the commercial data. Introducing PS within an IM framework involves
118 conditioning results upon structural hypotheses and then increases the computational
119 cost. We therefore assess how perform a more parsimonious model that would ignore PS.
120 Last, in addition to the PS, the fishing locations can be controlled by other factors
121 independent from the species distribution (e.g. logistical constraints – see Girardin *et al.*,
122 2017; Ducharme-Barth *et al.*, 2022). We therefore assess how such process blurring strict
123 PS may affect the quality of inferences.

124 We demonstrate the flexibility of the approach by fitting the model to three different
125 important European demersal fishery resources in the Bay of Biscay: common sole (*Solea*
126 *solea*, Linnaeus, 1758), hake (*Merluccius merluccius*, Linnaeus, 1758) and squids
127 (*Loliginidae* family). With these contrasted examples, we illustrate the capacity of the

128 framework to handle multiple commercial fleets with potentially distinct PS intensities and
129 different fishing behaviors.

130 **2 MATERIAL AND METHODS**

131 **2.1 Spatial integrated model**

132 Below we provide the core elements of the modelling approach. Additional details are
133 provided in supplementary material (SM) 1. The model is structured in four layers:
134 observations (here commercial and scientific CPUE in weight per unit of effort), sampling
135 process, latent field (here fish biomass relative density) and parameters (Figure 1 - all
136 notations are available in SM 1.1, Table S1). Sampling process is usually ignored in
137 hierarchical models as it is mostly considered independent of the quantity of interest, and
138 then has no consequence on the estimation procedure (Diggle *et al.*, 2010). Here, the
139 spatial distribution of commercial fishing is explicitly modelled as a non-homogenous
140 Poisson point process whose intensity may depend on the biomass field and contributes
141 to the likelihood. The observation processes of scientific and commercial data are
142 conditional upon the biomass latent field and the sampled locations.

143 All processes are considered to occur in a discrete fine grid (see for instance SM 2.1,
144 Figure S2.1 or SM 3.1, Figure S3.1). We assume the density of the point process is
145 piecewise constant in each cell grid which brings simplification in the expression of the
146 likelihood of the point process (Diggle, 2013 - see SM 1.2). The time component is omitted
147 and both commercial and scientific data are assumed to occur at the same time step.

148 The IM is designed to assimilate the scientific data of several surveys and/or the
149 commercial data of several fleets. In the following, the subscript j refers to the different
150 data sources either scientific or commercial. For instance, in a model with one scientific
151 survey and two commercial fleets, j will take the values $j = 1,2,3$, with $j = 1$ for the
152 scientific data and $j = 2,3$ for the two commercial fleets.

153 **2.1.1 Latent field of relative biomass**

154 The fish biomass relative density S (eq. (1) – (2)) is modeled through a latent log Gaussian
155 spatial field defined on the same discrete spatial domain as the point process. The mean
156 of the Gaussian field depends on environmental covariates through a log link where the
157 linear predictor combines an intercept α_S , the linear effect of environmental covariates
158 $\Gamma_S(x)$ (effects captured by the corresponding fixed parameters β_S representing the
159 species-habitat relationship). The remaining spatial variation is accounted for through a
160 zero-mean Gaussian random field (GRF) denoted $\delta(x)$ parameterized with a Matérn
161 correlation function $M(x, x'; \kappa, \phi)$, characterized by the shape κ and the scale ϕ (Cressie,
162 1993; Gelfand *et al.*, 2010; Lindgren *et al.*, 2011 and Banerjee *et al.*, (2014)). The shape
163 can be expressed in term of range $\rho = \frac{\sqrt{8}}{\kappa}$ where ρ is the distance for which the correlation
164 between points is near 0.1.

165
$$\log(S(x)) = \alpha_S + \Gamma_S(x)^T \cdot \beta_S + \delta(x) \quad (1)$$

166
$$\delta(x) \sim GRF(0, M(x, x'; \kappa, \phi)) \quad (2)$$

167 **2.1.2 Sampling process**

168 Recent literature has emphasized the complexity of the targeting behavior processes
169 (Salas and Gaertner, 2004; Haynie *et al.*, 2009; Abbott *et al.*, 2015; Girardin *et al.*, 2017;
170 Hintzen, 2021). In this paper, we did not attempt to model explicitly all those processes
171 and opted for a simplified representation where the spatial targeting directly depends on
172 the biomass field S and on an additional spatially structured random term.

173 Let us denote $X_{com,j}$ the spatial point process where commercial vessels of fleet j are
174 identified as fishing. In the following, all vessels in the same commercial fleet are assumed
175 to have homogeneous behaviors. Following Diggle *et al.* (2010), the set of fishing

176 locations are modeled conditionally on S , as a non-homogeneous Poisson point process
177 with piecewise constant intensity $\lambda_j(x)$ (eq. (3) - (4)).

$$178 \quad X_{comj} \sim \mathcal{JPP}(\lambda_j(x)) \quad (3)$$

$$179 \quad \log(\lambda_j(x)) = \alpha_{xj} + b_j \cdot \log(S(x)) + \eta_j(x) \quad (4)$$

180 For any fleet j , intensity $\lambda_j(\cdot)$ of the Poisson point process is modeled as a log-linear
181 combination of the logarithm of the relative biomass $S(\cdot)$ scaled by a parameter b_j , and a
182 residual spatial effect $\eta_j(\cdot)$ with the same structure as $\delta(\cdot)$ but with specific parameters κ
183 and ϕ . All parameters α_{xj} , b_j and the spatial random effect $\eta_j(x)$ are specific to each fleet.
184 The parameter b_j quantifies the strength of PS by scaling the relationship between the
185 local value of the resource field and the local fishing intensity.

186 Fishing locations potentially depend on many other factors than fish distribution such as
187 distance to harbor, logistical constraints, management regulations - spatial closures,
188 quotas – or fishing habits/tradition (Salas and Gaertner, 2004; Haynie *et al.*, 2009; Girardin
189 *et al.*, 2017). The spatial random effect $\eta_j(\cdot)$ is needed to capture any remaining additional
190 effect not captured by the dependence to $S(\cdot)$.

191 In that sense, a zero value for b_j indicates that the choice of the sampling locations does
192 not depend on the fish biomass relative density but only on the spatial random effect.

193 In addition to b_j , a dimensionless spatial metric was developed to quantify the strength of
194 PS (SM 1.3).

195 **2.1.3 Observation process**

196 Both scientific and commercial observations are considered as proportional to the
197 underlying biomass through a zero-inflated observation process. In our applications,

198 observations are expressed as CPUE (in weights / unit effort), with high proportion of
 199 zeros (zeros represent on average 30% of the commercial data and 10 to 50% of scientific
 200 data).

201 Observations are modelled through a zero-inflated lognormal model conditionally on
 202 biomass $S(x)$ in cell x (eq. (5-6)). The model is derived from Thorson *et al.* (2016) or
 203 Thorson (2018). We assume that the expected catch $\mu_j(x)$ for any fleet/data source j in
 204 the cell x depends on the latent field value $S(x)$ and a catchability coefficient q_j (eq. (5)).
 205 A zero catch ($y = 0$) is modeled as a Bernoulli random variable with parameter $\exp(-e^{\xi_j} \cdot$
 206 $\mu_j(x))$, where ξ_j is the parameter controlling the intensity of zeros relatively to the
 207 expected catch (eq. (6)). Then, $\mu_j(x)$ being fixed, the higher (resp., the lower) ξ_j , the lower
 208 (resp. the higher) the probability of obtaining a zero-catch.

209 The distribution of a positive catch $y > 0$ at a given x is defined as the combination of the
 210 probability of obtaining a non-zero catch ($1 - \exp(-e^{\xi_j} \cdot \mu_j(x))$) times a positive
 211 continuous distribution L (here a lognormal distribution) with expected value
 212 $\frac{\mu_j(x)}{(1 - \exp(-e^{\xi_j} \cdot \mu_j(x)))}$ and standard deviation σ_j . This formulation allows to represent the zero
 213 catch while assuring that the expected catch still equals $\mu_j(x)$.

$$214 \quad \mu_j(x) = q_j \cdot S(x) \quad (5)$$

$$215 \quad P(Y = y|x, S(x)) =$$

$$216 \quad \begin{cases} \exp(-e^{\xi_j} \cdot \mu_j(x)) & \text{if } y = 0 \\ \left((1 - \exp(-e^{\xi_j} \cdot \mu_j(x))) \cdot L \left(y, \frac{\mu_j(x)}{(1 - \exp(-e^{\xi_j} \cdot \mu_j(x)))}, \sigma_j^2 \right) \right) & \text{if } y > 0 \end{cases} \quad (6)$$

217 Per se, catchability q_j are not identifiable as there is no information in the model to
218 estimate the absolute scale of S . Commercial catches and/or scientific surveys will be only
219 informative about fish biomass relative density and additional information must be
220 provided to ensure statistical identifiability. If only one data type feeds the model (only
221 scientific or commercial data), relative catchability is fixed to 1 and the spatial random field
222 values is in the same scale as the data. If two data types (or more) are used to feed the
223 model, one of the relative catchability (denoted q_{ref}) has to be fixed, the other ones being
224 estimated relatively to the first one through a scaling factor k_j (eq. (7)).

$$225 \quad q_j = k_j * q_{ref} \quad (7)$$

226 As it is illustrated further in the simulation-estimation study (see section 3.1.1), the choice
227 of the reference level can have important consequences on the precision of estimation.

228 **2.1.4 Maximum likelihood estimation**

229 The estimation of the model is performed with TMB (Template Model Builder - Kristensen
230 *et al.* (2016)) and the spatial random effects are estimated through the SPDE approach
231 (Lindgren *et al.*, 2011) within the R software (R Core Team, 2020). More details on
232 estimation are available in the supplementary material (SM 1.4).

233 **2.1.5 Integrated model validation**

234 A key issue with IM is whether the different data sources provide consistent or conflicting
235 information (Saunders *et al.*, 2019; Zipkin *et al.*, 2019; Peterson *et al.*, 2021). In our
236 framework, the key question is whether integrating commercial data in addition to scientific
237 data will complement or will disrupt the inferences obtained from the scientific data,
238 considered as a reference source of information. To address this issue, we propose a
239 validation procedure based on the consistency check initially developed by Rufener *et al.*

240 (2021) and designed to check whether estimates obtained from the IM are consistent with
241 those obtained from the model fitted to scientific data only. The procedure would reject
242 consistency if the parameters estimates from the IM fall outside the 95% confidence region
243 of parameters estimates from scientific data only (see SM 1.5 for more details on the
244 procedure). This validation step is applied to both simulations and case studies.

245 **2.2 Simulation-estimation experiments**

246 We conducted simulation-estimation experiments to assess the performance of the
247 method for different data/model configurations (Table 1). Additional practical details on
248 the simulations are provided in SM 2. For all scenarios, simulations of data, covariates
249 and GRF were parameterized to tailor the case studies described hereafter. All scenarios
250 and configurations are repeated 100 times so as to capture the variability between
251 replicates.

252 Simulation-estimation experiments were specifically designed to address four questions
253 detailed below. In all cases, commercial data were simulated with various levels of PS
254 ($b = 0$ for uniform sampling, $b = 1$ for moderate PS, $b = 3$ for strong PS) to assess the
255 effect of PS on model's performance (Figure 2).

256 *(Q1) How do each data source contribute to inferences?*

257 In real case study, commercial data sample size may be far superior to scientific data
258 (specifically when using landings data) which might result in commercial data that
259 dominate inferences. To assess how the balance between the scientific and commercial
260 sample sizes drives the relative contribution of each data source, simulations were
261 conducted with few scientific samples (50 each) with increasing commercial samples
262 (50=small, 400=medium and 3000=large), and with a large commercial sample size

263 (3000) with increasing scientific sample size (50=small, 400=medium, 3000=large). No
264 scenario with more scientific samples than commercial samples is presented here as it is
265 a very unlikely configuration when using logbook catch data.

266 For each combination of commercial and scientific sample size, we fitted four different
267 models: a model fitted to scientific data only, a model fitted to commercial data only, and
268 two IM fitted to both commercial and scientific data, one with the scientific data used as
269 reference level and another one using the commercial data as reference level (Cf. eq. (7)).

270 For questions Q2, Q3 and Q4, all simulations were conducted using $n_{scientific} = 50$ and
271 $n_{commercial} = 3000$ to tailor the case studies. Commercial data are used as the reference
272 for catchability in the IM.

273 *(Q2) How does a partial coverage of the study area by the commercial data affect*
274 *the quality of the estimation?*

275 While scientific surveys are supposed to cover the full population distribution area, partial
276 coverage of the area by commercial fishing boats may arise from different sources like
277 spatial management closures (e.g. box closure) or too expensive travels from the coast.
278 To assess how a partial coverage by commercial data can affect estimates, we simulated
279 data with the commercial sampling intensity arbitrarily fixed to 0 in a fixed 9x9 box (15%
280 of the domain) while some biomass and some scientific samples are still simulated in this
281 area. We compared estimates of the biomass in the entire area with those obtained with
282 commercial data available on the whole domain.

283 *(Q3) What is the cost of ignoring PS in estimation when sampling is preferential?*

284 Modelling preferential sampling involves conditioning results upon a specified structural
285 assumption about sampling as well as increased computational cost. Here, we assess

286 how much ignoring PS would affect the quality of inferences when sampling is actually
287 preferential. We voluntarily introduce misspecification between the model used for
288 simulating the data (with various levels of PS intensity) and the one used in the estimation
289 procedure (b is alternatively estimated or arbitrarily fixed at 0).

290 *(Q4) How does the estimation perform when additional processes other than PS*
291 *drive the fishing locations?*

292 Fishing locations potentially depend on many other factors independent from the species
293 distribution (Salas and Gaertner, 2004; Haynie *et al.*, 2009; Girardin *et al.*, 2017). To
294 assess how such process blurring strict PS may affect the quality of inferences, we
295 simulate data with a sampling intensity that depends on both the biomass distribution (PS)
296 and an additional spatial random terms $\eta_f(\cdot)$ independent from the biomass distribution
297 (eq. (4); see Table 1 for more details on $\eta_f(\cdot)$ parameterization), and compare the
298 inferences obtained from a data set simulated with strict PS ($\eta_f(\cdot) = 0$ on the full domain).
299 Note that for questions Q1, Q2 and Q3, the random effect η was fixed to 0 in simulations
300 (but it is still estimated in the estimation model), so that the sampling process only
301 depends on the distribution of biomass.

302 **2.2.1 Performance metrics**

303 The performance of the estimation method was assessed using different metrics on key
304 model parameters.

305 The quality of the total biomass estimation (the sum over all grid cells, $B = \sum_x S(x)$) was
306 explored through the relative bias $\frac{(B-\hat{B})}{B}$, that quantifies how much the total biomass is over
307 or under-estimated.

308 The quality of the estimation of the parameter b is assessed through the relative bias
309 defined as $\frac{b-\hat{b}}{b}$ (except for $b = 0$, where only the absolute bias is considered). We also
310 assessed the relative bias of the species-habitat relationship estimate $\hat{\beta}_S$ and range
311 parameter ρ as these parameters are meaningful for understanding species distribution.
312 The precision of the spatial predictions was studied with the mean squared prediction error
313 between the simulated and the estimated latent field values $\frac{1}{n} \sum_x (S(x) - \widehat{S}(x))^2$ (MSPE –
314 n stands for the number of grid cells).

315 **2.3 Case studies**

316 We applied the approach on three case studies of demersal fisheries in the Bay of Biscay:
317 the common sole (*Solea solea*, Linnaeus, 1758), the hake (*Merluccius merluccius*,
318 Linnaeus, 1758) and the squids (Loliginidae family). These case studies were selected
319 because they emphasize different intensities of preferential sampling. Further details on
320 case studies and data are provided in SM 3.

321 To compare models on the same spatial domain for the three species, we limited the
322 analysis to scientific and commercial data available on the Bay of Biscay only (SM 3.1,
323 Figure S3.1 for the spatial grids). Besides, to get some replicates of the analysis, we
324 applied the approach on 2 years for each case study (2017 and 2018 for common sole –
325 2014 and 2015 for hake and squid). To keep it synthetic, only the data and the results of
326 the models for hake in 2014, sole in 2017 and squids in 2015 are presented in this
327 manuscript. The related IM estimating PS passed the consistency check for both statistical
328 tests and they allow to clearly illustrate the effect of PS on model outputs.

329 **2.3.1 Survey data**

330 Scientific data (CPUE, in kg/hour) were derived from the Orhago survey for common sole
331 and EVHOE survey for hake and squids (ICES, 2020a). The sampling density (number of
332 data points / km²) of those two surveys revealed representative of the sampling density of
333 the main European trawl surveys from the DATRAS database (see SM 3.2). In
334 comparison, commercial data used in the case studies are denser by 2 orders of
335 magnitude. Scientific data was aligned on commercial data by filtering only individuals
336 above the minimum landing size when available (24 cm for sole, 27 cm for hake - ICES,
337 2020).

338 **2.3.2 Commercial data**

339 For each species, we filtered commercial data for 'bottom trawlers' as they cover a wide
340 part of the study area (Figure 3) and provide easy to compute and reliable CPUE.
341 Commercial data were standardized by the fishing effort in (kg/hour). For hake and sole,
342 we filtered the métier targeting demersal fish (called OTB_DEF) and for squids, the métier
343 targeting cephalopods (called OTB_CEP).

344 The orders of magnitude of commercial sample size is much higher than for scientific data.
345 For hake (i.e. OTB_DEF), there are 6852 commercial samples in 2014 and 5000 in 2015.
346 For squid (i.e. OTB_CEP), there are 7486 commercial samples in 2014 and 9611 in 2015.
347 This should be compared with the 86 EVHOE samples for both years. For sole (i.e.
348 OTB_DEF), there are 2401 samples in 2017 and 3325 in 2018 compared with the 49
349 Orhago samples for both years.

350 **2.3.3 Habitat covariates**

351 Two covariates classically used to describe benthic species distribution were selected:
352 depth and sediment type (Le Pape *et al.*, 2003; Witman and Roy, 2009; Rochette *et al.*,
353 2010). Depth was separated into several categories and was considered (as sediment)
354 as a categorical variable (SM 3.7, 3.8).

355 **2.3.4 Model configurations**

356 As for the simulation-estimation experiments, the models were fitted under different
357 configurations. To assess the information brought by each dataset, we compared the
358 model fitted to scientific data only, to commercial data only and to both scientific and
359 commercial data. To assess the effect of PS on model outputs, we compared IM
360 accounting for PS (b is estimated) with IM where PS is ignored (b is fixed to 0).

361 For the sole case study, we compared results obtained from the IM by considering one
362 homogeneous or two distinct fleets with specific catchability and targeting parameters.
363 Note that splitting one fleet in 2 distinct fleets is performed through a PCA coupled with a
364 HCPC analysis on vessels characteristics data derived from both logbooks and VMS data.
365 All the clustering analysis is described in SM 3.9.

366 **2.3.5 Model evaluation**

367 Uncertainty of the predictions are quantified through the coefficient of variation and all
368 estimates (e.g. fixed parameters, total biomass) are represented with related 95 %
369 confidence intervals. We assess the consistency of IM through the statistical tests
370 described in section 2.1.5 and in SM 1.5. Finally, the different IM are compared through a
371 5-fold cross validation, and model performance was quantified based on two metrics: the
372 $MSPE_{fit}$ that measures goodness of fit (MSPE – mean squared prediction error), and the

373 *PCV* that measures predictive capacity (see SM 3.10 for more details on the metrics and
374 guidelines for interpretation).

375 **3 RESULTS**

376 **3.1 Simulations**

377 We summarize the main results of the simulation-estimation experiments below.

378 Additional results are provided in SM 4.

379 ***3.1.1 Contribution of each data source in the integrated model***

380 Models fitted on scientific data only provide systematically unbiased estimates of total
381 biomass (the mean bias is close to 0 for all sample size - Figure 4, 1st row), and the
382 variance of estimations logically decreases with scientific sample size. Note that the
383 species-habitat relationship estimates $\hat{\beta}_S$ are also unbiased (see SM 4.1).

384 Overall, inferences from IM revealed consistent with those obtained from scientific data
385 only (SM 4.2.1). Even when the commercial sample size is high and the scientific sample
386 size is low, only 3% of the p-values fall below the 0.05 threshold for the fixed effect test
387 (the test wrongly rejects consistency). For the random effect test, the results are more
388 balanced as 10% of the p-values fall below the 0.05 threshold when data size are very
389 unbalanced (low scientific sample – high commercial sample).

390 In almost all configurations, IM provide unbiased and more precise estimates for total
391 biomass and spatial biomass predictions compared to the model fitted to scientific data
392 only (Figure 4). As expected, the higher the commercial and the scientific sample size,
393 the more accurate the spatial predictions, the PS parameter b and total biomass
394 estimates. Estimates of b are unbiased in most cases except when commercial sample
395 size is low and PS is strong (Figure 4, 2nd row).

396 As expected, the contribution of each data sources in the IM directly depends on the
397 balance in the sample size. When sample size is balanced between the data sources,
398 then integrating the two data sources in the model systematically improves the inferences
399 with regards to situations where only one data source is analyzed. For instance, for large
400 commercial and scientific sample size (com.L_sci.L) and no PS, the precision is 1.5 higher
401 (i.e. the MSPE is 1.5 lower) for the IM compared to single-data models (either scientific or
402 commercial - Figure 4, 3rd row, 1st column). However, when the sample sizes are
403 unbalanced, the data source with the higher sample size (here commercial data)
404 dominates inference and integrating another data source with a smaller sample size (here
405 scientific data) contributes to a much lesser extent to inference. See for instance the
406 situation where commercial sample size is large and scientific sample size is low
407 (com.L_sci.S - Figure 4, 3rd row, 1st column). In this case, the performances of the model
408 fitted to commercial data alone – with reference level fixed to commercial data - are very
409 close to those of the IM whatever the intensity of PS.

410 Interestingly, the higher the intensity of PS, the higher the benefits of integrating
411 commercial data in the model (Figure 4, 3rd row); for instance, when both datasets have
412 large sample sizes (com.L_sci.L), increasing PS reduces error predictions (i.e. increases
413 accuracy) by 2 each time (i.e. for $b = 0$, $E(MSPE) = 20$; for $b = 1$, $E(MSPE) = 10$; for $b =$
414 3 , $E(MSPE) = 5$).

415 Still, the simulations also reveal some limits in the inferences. First, the range parameter
416 might be poorly estimated and slightly biased when the sample size is low while being
417 better estimated when increasing the sample size or integrating additional data in the
418 analysis (see SM 4.3).

419 Also, in unbalanced cases the accuracy of total biomass estimates from the IM revealed
420 highly sensitive to the choice of the reference level (Figure 4, 1st row). When the
421 commercial sample size far exceeds the scientific sample size, setting the reference level
422 to the commercial data produces more precise estimates than setting the reference level
423 to scientific data. When defining scientific data as reference level, the intercept of the
424 latent field of relative biomass is estimated from the few scientific samples and resulting
425 estimates are less precise than when defining the reference level with a more numerous
426 data source (here commercial data). This is also true - to a lesser extent - for spatial
427 predictions (Figure 4, 3rd row).

428 In the following, only the case where commercial samples exceed scientific samples and
429 the reference level is fixed with commercial data is explored further as it is the closest to
430 the case studies configuration (Table 1).

431 **3.1.2 Impact of a partial coverage of the study area by the commercial data**

432 When commercial data only partially cover the distribution area, commercial data still
433 provide valuable information to predict biomass spatial distribution whatever the PS
434 intensity (Figure 5, 2nd column). When sampling is not preferential (data simulated with
435 $b = 0$), a partial coverage of the distribution area produces on average 1.5 less precise
436 spatial predictions but estimates remain unbiased (Figure 5, 3rd row, comparing 1st and
437 2nd column). When sampling is preferential (either moderate or high), biomass estimates
438 are slightly underestimated. Integrating scientific data in the analysis does not correct this
439 bias.

440 Finally, all model configurations allow for unbiased and precise estimation of the species-
441 habitat parameters $\hat{\beta}_S$ whether or not there is a partial coverage of the domain (see SM
442 4.1) and overall almost all IM are consistent with scientific-based model (SM 4.2.2).

443 **3.1.3 How does ignoring PS impact inferences?**

444 As expected, the impact of ignoring PS in the estimation model is negligible when data is
445 simulated with no PS, and becomes more and more detrimental when the intensity of PS
446 increases in the truth (Figure 5, 3rd column). With no surprise, when data are generated
447 with no PS ($b = 0$), ignoring PS in the estimation procedure has no effect on the estimation
448 performance. When PS is moderate, total biomass estimates are 5 % overestimated ($b =$
449 1). In the case of strong PS ($b = 3$), ignoring PS in the estimation strongly deteriorates the
450 quality of inferences regarding total biomass estimates (Figure 5, 1st row, 3rd column).
451 Total biomass estimates are overestimated by 50% on average. However, the main spatial
452 patterns are well identified with or without consideration of PS, even though more precise
453 when accounting for PS (Figure 5, 3rd row, 1st column). SM 4.4 (Figure S4.4.1) presents
454 maps comparing a simulated biomass field and model predictions obtained by considering
455 or ignoring PS when $b = 3$. The areas with high biomass values (i.e. where commercial
456 sampling is dense) are well predicted by the models accounting for PS or not. The main
457 differences are localized in poorly sampled areas where biomass is low. Accounting for
458 PS in estimation allows to interpret the low sampling intensity areas as low-density areas,
459 and therefore to reduce the bias in those areas (SM 4.4, Figure S4.4.2).

460 Finally, from a computational point of view, accounting for PS on average multiplies by 4
461 the computational time (see SM 4.5).

462 **3.1.4 Effect of other spatially structured processes affecting fishing locations**

463 As expected, precision of estimates are deteriorated when fishing locations actually
464 depend upon a combination of biomass distribution (PS) and other mechanisms (here
465 captured by a spatially structured random term - Figure 5, 4th column). In this case, the
466 IM still provides valuable inferences on fish distribution, fish total biomass and estimates
467 of b , although estimations are less accurate than the base case. For instance, MSPE are
468 5 times lower when nothing else than PS affects sampling locations compared with a case
469 where sampling locations depend on both PS and other independent spatial processes
470 (Figure 5, 3rd row, 1st and 4th column). But interestingly, the weight of scientific data
471 increases when the sampling distribution of commercial data is blurred by spatial
472 processes independent from biomass spatial distribution. MSPE and relative bias
473 provided by the IM are both 1.4 smaller compared to those obtained when the model is
474 fitted to commercial data only.

475 **3.2 Case studies**

476 Below we summarize the main results obtained from the application of the framework to
477 the three case studies. Additional results and maps are provided in SM 5.

478 **3.2.1 Contribution of each dataset to the inferences**

479 Almost all the case studies successfully passed the consistency test between the IM and
480 the model fitted to scientific data only (see SM 5.1). Still, models based on scientific data
481 provide different spatial predictions compared with the IM. Predictions for sole and squids
482 from the scientific-based model are mainly shaped by the covariate effects (Figure 6; for
483 further analysis see SM 5.2, SM 5.3 and SM 5.4). On the other hand, predictions from the
484 IM are mainly shaped by the spatial random effect as commercial data allow to better
485 capture the local spatial correlation structures.

486 Consistently with simulations, inferences from the IM are mainly driven by the commercial
487 data (Figure 6). This logically arise from the much higher sample size of commercial data
488 compared with scientific data, combined with the good coverage of commercial data in
489 high-density areas (Figure 3). As commercial data is denser than scientific data, they will
490 better capture local spatial correlation structures than scientific data. SM 5.5 provide some
491 additional analysis of the information brought by commercial data in the IM.

492 In this configuration, scientific data bring information to model predictions in areas poorly
493 covered by the commercial data (SM 5.6 - e.g. for squids, the offshore predictions are
494 downscaled by scientific data).

495 **3.2.2 Preferential sampling and other processes affecting fishing locations**

496 In this section and related SM (SM 5.7 to SM 5.10), we focus on results from the IM only.

497 For the three case studies, estimates of b are positive, suggesting sampling by fishermen
498 is preferential towards high biomass density areas. The hake case study has the lowest
499 PS parameter ($\hat{b} = 0.88$, $sd(\hat{b}) = 0.107$), followed by sole ($\hat{b} = 2.4$, $sd(\hat{b}) = 0.046$) and
500 squids ($\hat{b} = 3.5$, $sd(\hat{b}) = 0.025$). For more intuition concerning the strength of PS and
501 how it varies in space, refer to SM 5.7. In all case studies, the spatial random term η in
502 the sampling process turned to be spatially structured (SM 5.8) and captures 25% to 97%
503 of the spatial variability of fishing locations (SM 5.9). This highlights the importance of
504 other spatial mechanisms in the choice of fishing locations compared to strict PS towards
505 biomass distribution.

506 Consistently with simulations, the higher the PS intensity, the higher the differences
507 between inferences obtained with and without considering PS. When comparing biomass
508 field values (Figure 7, left column), ignoring PS increases predictions in poorly sampled

509 areas (all red areas – compare with Figure 3). This effect is particularly marked for the
510 squid case study where the relative difference is the strongest in the offshore areas.
511 However, considering PS or not has relatively little effect in areas where sampling is
512 spatially denser (all white areas). Ignoring PS affects total biomass indices estimates and
513 the relative difference between biomass estimates with or without PS increases with the
514 value of b estimates (Figure 7, right column).

515 When the estimated PS intensity is high (i.e. in the case of squids) accounting for PS can
516 improve model goodness-of-fit and predictive capacity (SM 5.10).

517 **3.2.3 Benefits of considering different fleets in the estimation model**

518 Based on the sole case study, we demonstrate the capacity of the model to integrate
519 multiple commercial fishing fleets, each with specific parameters (catchability and
520 targeting). In the sole case studies, considering two different fleets in the IM (instead of
521 one homogeneous) improves goodness-of-fit towards scientific data (SM 5.11, y-axis) and
522 modifies spatial predictions (SM 5.12).

523

524 **4 DISCUSSION**

525 **Main findings**

526 Combining multiple sources of data to build more informative spatio-temporal models for
527 fish distribution is a major challenge in fishery ecology. Commercial catch per unit effort
528 data have long been recognized as a valuable source of information eventually highly
529 complementary to scientific survey data. But the complexity of the mechanisms driving
530 the way fishermen sample in space and time make the combination of scientific and
531 commercial data challenging.

532 In this paper, we provide a hierarchical framework to integrate scientific surveys and
533 commercial catch declaration data to infer species distribution while considering the effect
534 of PS on fishing points distribution. The new model allows for exploring and questioning
535 the challenges raised by such integration. The benefit but also the limits of the new
536 approach were evaluated using simulations and through the application of the model to
537 three contrasted demersal case studies (sole, hake and squids) of the Biscay Bay fishery.

538 Both simulations and case studies demonstrate that ignoring PS in the inference may be
539 highly detrimental when the intensity of PS is strong. The present framework can serve
540 as a tool to assess the benefit of including PS in analysis, depending on the intensity of
541 PS but also on the modelling objectives. As already shown in previous studies (Conn *et*
542 *al.*, 2017; Pennino *et al.*, 2019), when PS actually occurs in commercial catches, ignoring
543 this process may bias inferences on total biomass estimates. Even if ignoring PS may not
544 hamper the capacity to detect areas of high biomass, the biomass in low-density areas
545 may be overestimated. Therefore, if the objective is to compute biomass indices integrated
546 over a large area, then it might be worth accounting for PS to avoid biased results. By

547 contrast, if the objective is to identify hotspots, the benefits of considering PS may be low
548 with regard to the additional computational time it requires.

549 The three case studies illustrated the potential of the model to handle the variability of PS
550 behavior among species and fleets. Low PS was revealed for hake, while a moderate and
551 strong PS was revealed for sole and squids, respectively, which is consistent with the
552 expert knowledge on the behavior of those bottom trawls fleets (Y. Vermard, *com. pers.*).
553 Results also demonstrate the capacity of the framework to integrate commercial catch
554 data from multiple fleets, and the benefits for the quality of inferences when those fleets
555 have different features such as distinct catchabilities or targeting behaviors. For the sole
556 case study, this approach proves useful to distinguish two segments in the bottom trawl
557 fleet, which improved model outputs. This framework could be extended to more than two
558 fleets and combined with other studies analyzing fleets structure (Pelletier and Ferraris,
559 2000; Ferraris, 2002; Stephens and MacCall, 2004; Deporte *et al.*, 2012; Winker *et al.*,
560 2013; Okamura *et al.*, 2018).

561 **Challenges in modelling PS**

562 Still, modelling the spatial distribution of commercial fishing locations remains highly
563 challenging (Hintzen, 2021; Hintzen *et al.*, 2021). Our framework is shaped to integrate
564 data from homogeneous fishing fleets supposed to share the same fishing behavior, which
565 simplifies the modelling of the non-uniform spatial intensity of fishing for each fleet. We
566 propose a parsimonious model where the dependence of the sampling intensity to the
567 biomass is supposed to be linear in the log scale. This is a strong hypothesis and
568 departure from this hypothesis may obviously exist in the truth. For instance, the intensity
569 of PS could vary in space such as in Conn *et al.* (2017) who considered that the degree

570 of PS could change across the landscape. On the other hand, however, the log-log linear
571 assumption is easy to implement in other software including the VAST R package used
572 for operational assessments in some management regions (Thorson *et al.*, 2019).

573 Of course, many other factors may drive the spatial intensity of fishing, and those were
574 simply captured in our model through an additional spatial random term. For instance,
575 fishers' behavior may depend on prior knowledge of fish spatial distribution, on information
576 sharing within fishing cooperatives, on expected distribution of bycatch species, or
577 logistical constraints (e.g., transit costs) (Salas and Gaertner, 2004; Haynie *et al.*, 2009;
578 Girardin *et al.*, 2017). Targeting behavior may also be directed toward an assemblage of
579 species rather than toward a single species (Bourdaud *et al.*, 2019).

580 The random effect should be able to capture additional variations whenever the departure
581 from a continuous Gaussian random field is not too high. If not, for instance in the case of
582 fishery closures where fishing activity suddenly drops to very low levels (as explored in
583 simulation-estimation), the model may produce biased estimates due to model
584 misspecification. We did not detect such misspecification in our case study, but we
585 recommend that future analyses based on fishery-dependent data present a log-log plot
586 between sampling intensity and predicted biomass density to diagnose strong departure
587 from model hypothesis.

588 Still, some non-spatial targeting has been reported from multi-species catch records
589 (Stephens and MacCall, 2004; Okamura *et al.*, 2018). Efforts to integrate these methods
590 into spatio-temporal models are underway (Thorson *et al.*, 2016), although these methods
591 have not previously been extended to jointly analyzing multi-species fishery and survey
592 data.

593 **Relative contribution of scientific and commercial data**

594 Our analysis exemplifies that a key issue in such integrated modelling exercise is to get a
595 sensible evaluation of the relative contribution of the different sources of data in
596 estimation. In particular, critical issues with IM are whether the different data sources
597 provide eventually highly unbalanced quantity of information (then the inferences are fully
598 dominated by one of the data sources; Fletcher *et al.*, 2019) and whether they provide
599 complementary or conflicting information to the final inferences (Saunders *et al.*, 2019;
600 Zipkin *et al.*, 2019; Peterson *et al.*, 2021).

601 We implemented a likelihood ratio-test (Rufener *et al.*, 2021) to check for model
602 consistency between the IM and the scientific-based model. In most cases, models
603 passed the consistency check successfully, although it was rejected in some cases. Some
604 further analysis should investigate in detail the reasons of these inconsistencies as they
605 could probably shed light on some new research avenues for model improvement. For
606 instance, some neglected vessel effect (*e.g.*, difference in catchability among vessels -
607 Thorson and Ward, 2014) or some too simplistic representation of the sampling and/or
608 the observation process of commercial data might partly explain these inconsistencies.

609 Simulations revealed that when scientific data and commercial data have balanced
610 sample size, they both contribute to inference and the IM will provide better biomass
611 predictions than models based on single-data set. As expected, when the sample size of
612 commercial data far exceeds scientific data, inference about spatial patterns is mainly
613 driven by the commercial data. In the three case studies, we used commercial data with
614 sample sizes that far exceed the scientific one. In that case, scientific data have relatively
615 limited weight in the final inference. Still, they bring valuable information in areas that are
616 not sampled by the commercial fishery. Also, scientific data remain a critical component

617 in the analysis as they provide some reference data through a standardized sampling plan
618 and a controlled protocol allowing then to assess for the IM consistency. It would be worth
619 applying our framework to other case study that may consist in more balanced data sets,
620 such as models seeking to combine scientific with onboard observer data (Rufener *et al.*
621 2021), or in pelagic fisheries where acoustic surveys can provide continuous observations
622 over the full domain.

623 Our results also point out the importance of setting the reference level for the catchability
624 coefficient with either the scientific or the commercial data. In particular, when the sample
625 size of the commercial data far exceeds the scientific survey, fixing the reference level
626 with scientific surveys generally results in higher imprecision, due to the lower sample
627 size. But still, in certain cases, the scientific data may provide absolute information on
628 biomass and fixing the catchability factor associated with the survey data can result in an
629 interpretable measure of index scale (Thorson *et al.*, 2021). Hence, the choice of the
630 reference level could be a matter of tradeoffs between precision of inferences and
631 interpretation of the results in terms of scale.

632 **The limits of reallocated catch data**

633 Probably one of the major limits of our approach is that the actual framework ignores the
634 uncertainty that arises from the procedure used to reallocate the catch declarations in
635 space. Obtaining the spatialized CPUE inputs used in the model requires pre-treatment
636 of the commercial catch declaration data to allocate declaration data to VMS positions
637 (Hintzen *et al.*, 2012). Raw data corresponds to fishing operations that are daily
638 aggregated and reported at coarse administrative spatial units (0.5° latitude by 1°
639 longitude rectangles). These declarations are then reallocated uniformly on all GPS

640 locations previously identified as fishing in the vessel path. This procedure has been
641 demonstrated to be robust while being a fast and a pragmatic approach for reallocating
642 landings to VMS pings (Gerritsen and Lordan, 2010; Murray *et al.*, 2013). However, it
643 implies strong hypotheses that may artificially increase or transform the information
644 provided by the data. Typically, the uniform reallocation of catch declarations on all GPS
645 positions identified as fishing may smooth the spatial signal, which could potentially
646 explain the lack of species-habitat relationship obtained from the IM. The effect of such
647 reallocation should be explored in further study to better understand its consequences on
648 model predictions/estimates and further model development should investigate how to
649 mitigate its consequences.

650 **Perspectives**

651 Our work raises some major challenges which all constitutes exciting tracks for future
652 research.

653 Data-weighting approaches could be explored further to better control the contribution of
654 the two sources of data and eventually assess if increasing scientific data weight could
655 improve model predictive capacity. Data-weighting methods intend to modify the relative
656 influence of the data sources by assigning or estimating a weight for each data source
657 (Francis, 2017; Punt, 2017; Wang and Maunder 2017; Punt *et al.*, 2020). Only very few
658 studies have already explored the potential for data weighting in the SDM context
659 (Fletcher *et al.*, 2019). Still, several questions regarding the weight specification remain
660 open or largely debated. For instance, how to rigorously fix/estimate/interpret the weight?
661 Also, when can we consider that a data-weighting approach is relevant or is it only a matter
662 of model misspecification? Some theoretical and modelling development could be highly

663 valuable to provide a generic and rigorous formalization for either data weighting or model
664 correction in the context of SDM (but see for instance the approach provided by Thorson
665 *et al.* (2017b) for composition data in the context of stock assessment models).

666 Another option would consist in developing an alternative observation model for the
667 commercial CPUE in order to better capture the uncertainty associated with the
668 reallocation procedure. As a general idea, an observation model could be developed to
669 explicitly represent that CPUE are available at the scale of the daily fishing activity (the
670 scale that corresponds to the catch declaration), rather than artificially reallocating
671 uniformly catch declarations on related VMS pings. Doing so, the quantity of information
672 provided by commercial data would be more representative of the information they really
673 contain.

674 Future work should also seek to better integrate the discrete-choice and econometric
675 analyses emphasizing the complexity of the processes related to the choice of fishing
676 locations. For instance, the sampling process could account for the pluri-specific nature
677 of fisheries (Bourdaud *et al.*, 2019) and additional factors other than fish distribution could
678 be included to explain the variability of sampling intensity in space and time (Salas and
679 Gaertner, 2004; Haynie *et al.*, 2009; Girardin *et al.*, 2017).

680 Finally, including a temporal dimension in the model and fitting a longer time series looks
681 a fruitful research avenue. Moving to spatio-temporal modelling that would consider
682 temporal autocorrelation in the spatial distribution may be methodologically challenging
683 (Cameletti *et al.*, 2013), but represents an exciting step towards a better understanding of
684 the seasonal spatial distribution of fish resources. Indeed, commercial data are often
685 available all along the year, when scientific surveys most often occur once or twice a year.
686 Combining scientific and commercial data within an integrated spatio-temporal framework

687 built at an infra-annual time step (e.g., season, month) would allow to complement the gap
688 of information to investigate fish spatio-temporal distribution at a finer temporal scale than
689 what is possible using scientific data only (Hilborn and Walters, 2013; Maureaud et al.,
690 2020). It would offer new opportunities to interpret seasonal patterns of distribution (Kai et
691 al., 2017), identify fish functional habitats such as spawning areas (Paradinas et al., 2015;
692 Delage and Le Pape, 2016), and provide the required knowledge for protecting those
693 habitats (Schmitten, 1999; Erisman et al., 2020).

694 **SUPPLEMENTARY MATERIAL**

695 All the supplementary material documents are available at the ICESJMS online version of
696 the manuscript. They provide additional information on the modelling framework (SM1),
697 material and methods for simulations (SM2) and case studies (SM3), results for
698 simulations (SM4) and case studies (SM5).

699 **ACKNOWLEDGMENTS**

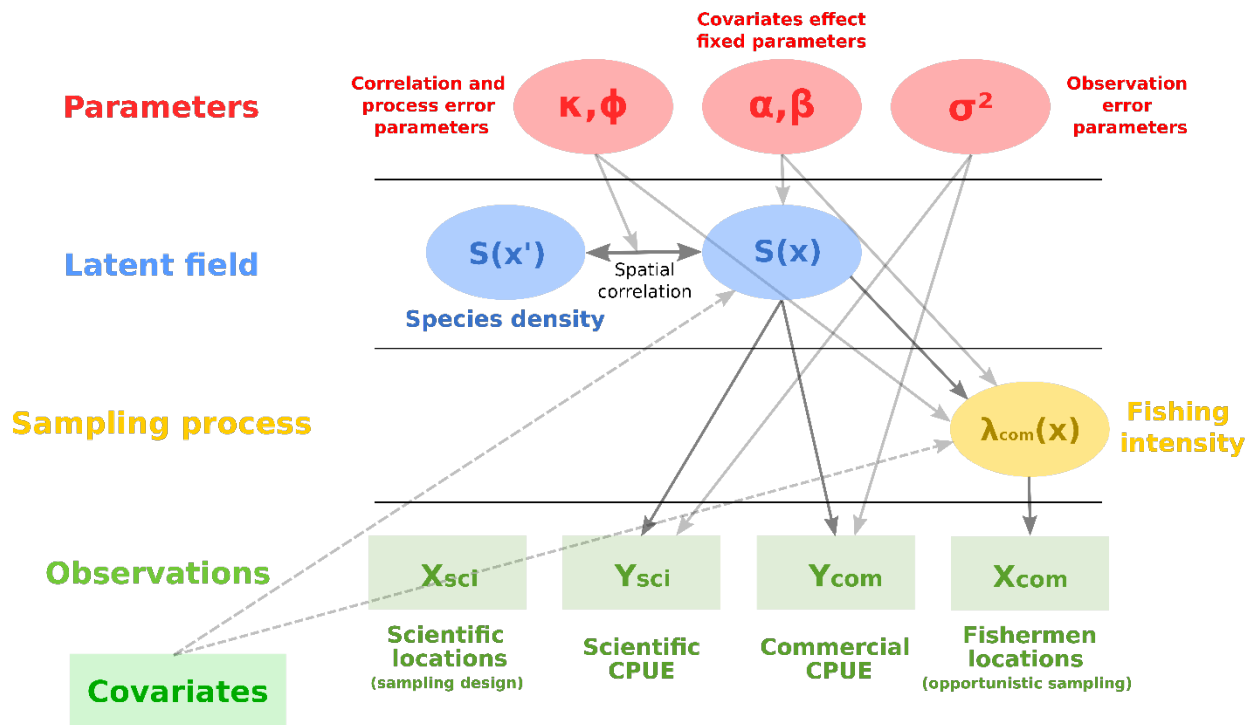
700 The authors are grateful to the Direction des pêches maritimes et de l'aquaculture (DPMA)
701 and Ifremer (Système d'Informations Halieutiques - SIH) who provided the aggregated
702 VMS data. The findings and conclusions of the present paper are those of the authors.

703 The authors thank David Eme who provided tidy environmental covariates data as well as
704 Kasper Kristensen, Jean-Baptiste Lecomte, Louise Day and Pierre-Yves Hervann for
705 their feedbacks and their highly valuable advice. The authors thanks also Maxime Olmos,
706 John Best, and two anonymous reviewers whose feedbacks greatly improved the
707 manuscript.

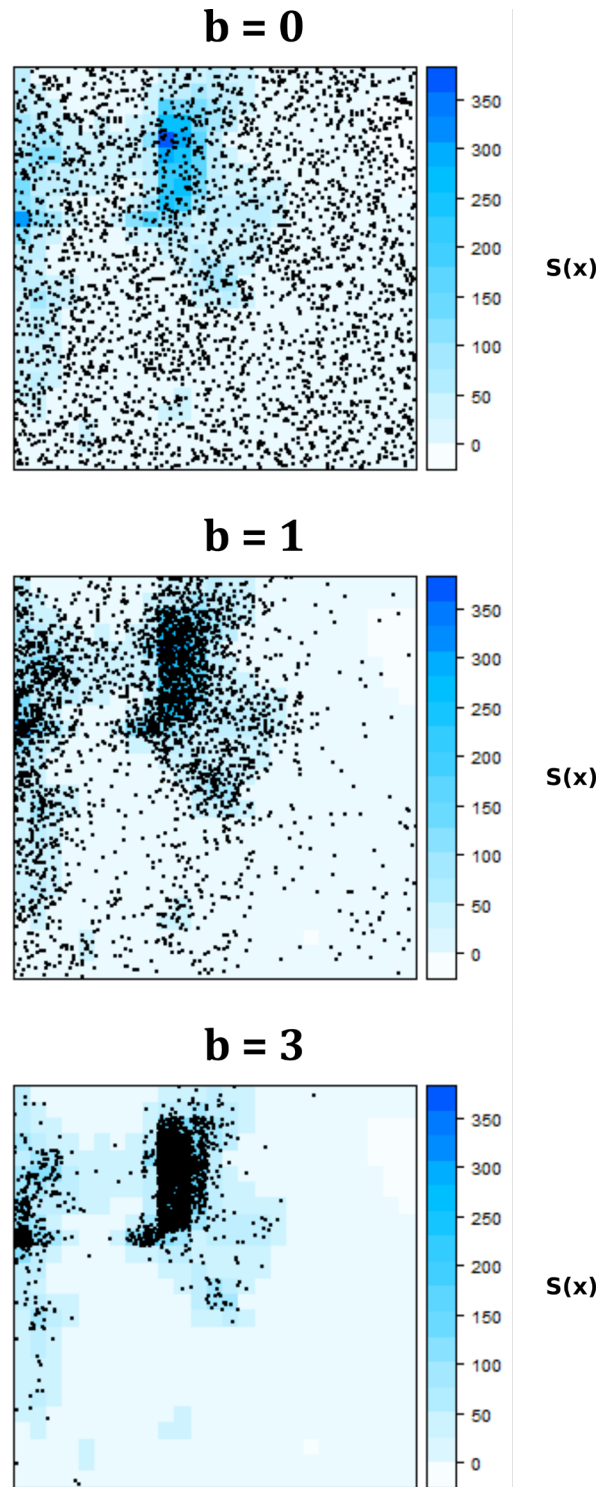
708 **DATA AVAILABILITY STATEMENT**

709 Survey data are available through the DATRAS portal ([https://www.ices.dk/data/data-](https://www.ices.dk/data/data-portals/Pages/DATRAS.aspx)
710 [portals/Pages/DATRAS.aspx](https://www.ices.dk/data/data-portals/Pages/DATRAS.aspx)) with the package 'icesDatras' ([https://cran.r-](https://cran.r-project.org/web/packages/icesDatras/index.html)
711 [project.org/web/packages/icesDatras/index.html](https://cran.r-project.org/web/packages/icesDatras/index.html)). Logbooks and VMS data are
712 confidential data and they are available on specific request to DPMA. Codes that support
713 the findings of this study are on gitlab and can be given access on request at the address:
714 baptiste.alglave@agrocampus-ouest.fr.

715



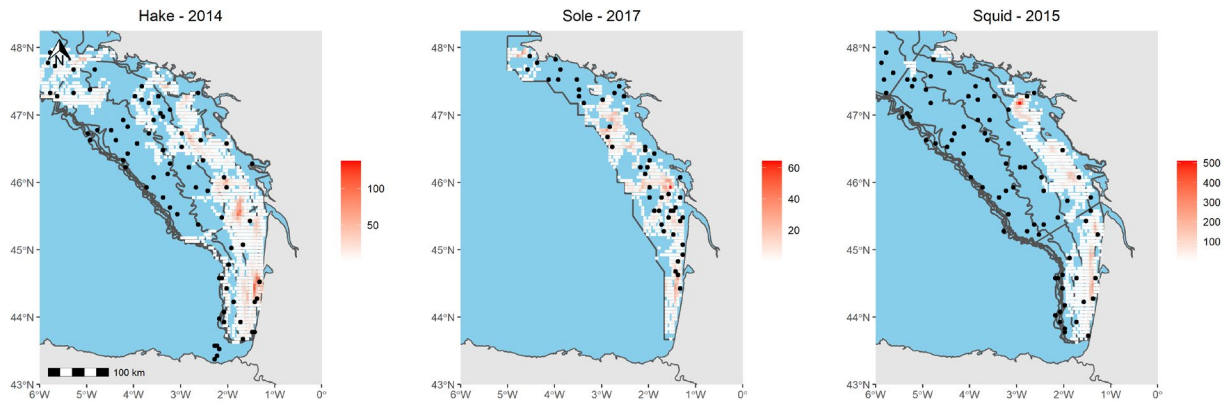
719 *Figure 1. Diagram of the spatial integrated model including preferential sampling for*
 720 *commercial data. Locations of scientific trawls do not contribute directly to the likelihood.*



722 *Figure 2. Maps of simulated commercial sampling points obtained for three values of*
 723 *preferential sampling ($b=0$, $b=1$, $b=3$). Blue scale: values of the simulated biomass field.*
 724 *Dots: fishing points. For $b = 0$, the targeting metric $T_j(x) = 1$ (SM 1.3). For $b = 1$,*

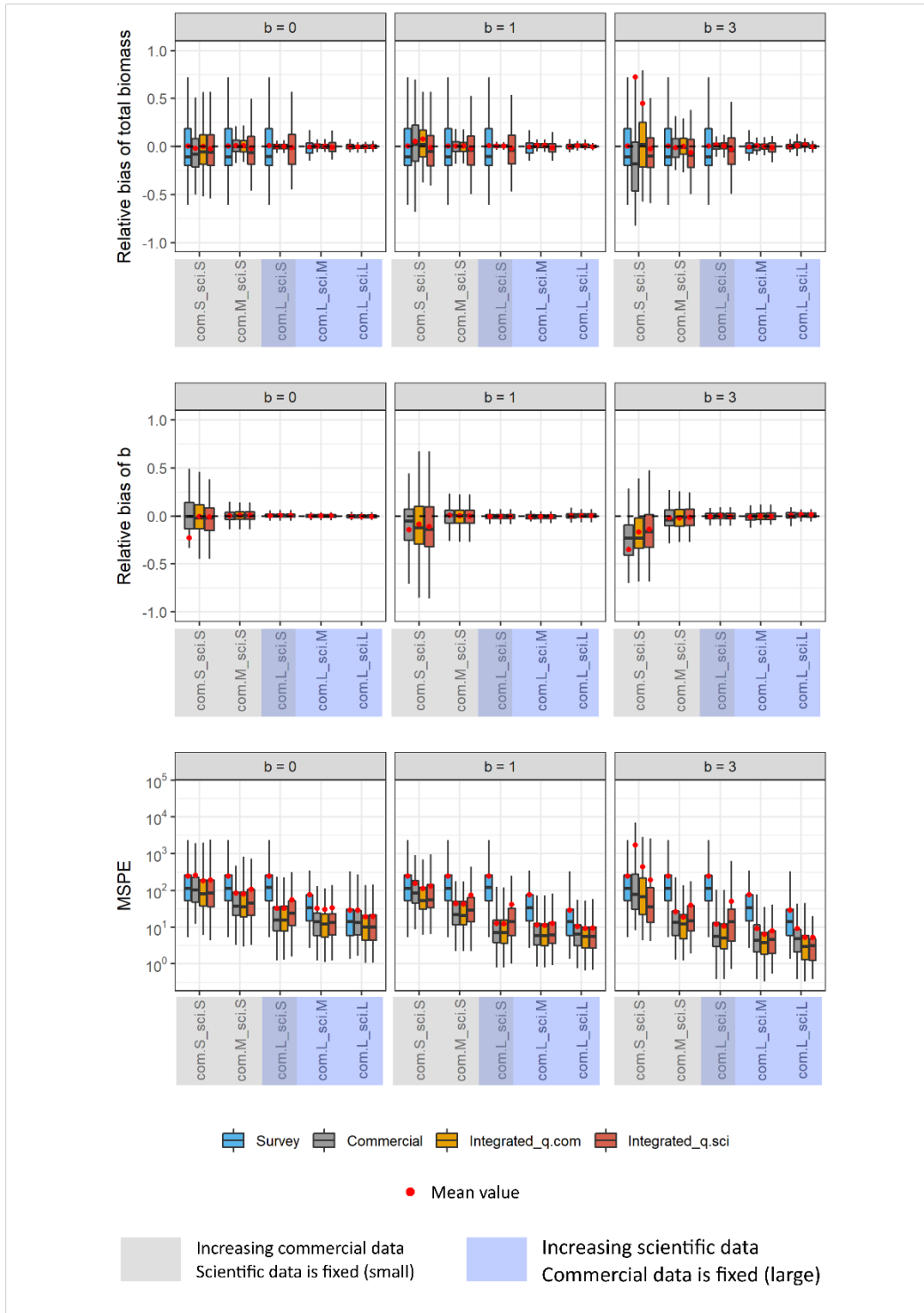
725 $\arg \max_x T_j(x) = 12, q_{50\%}\{T_j(x)\} = 0.4$. For $b = 3, \arg \max_x T_j(x) = 80, q_{50\%}\{T(x)\} =$
726 $0.002.$

727



729

730 *Figure 3. Map of scientific samples (black dot) and commercial sampling distribution*
 731 *(red color scale – unit: fishing hours). Note that all scientific hauls last around 30*
 732 *minutes. Black lines - limits of the spatial domains covered by the scientific survey*
 733 *(Orhago and EVHOE) that delineate the study area. Left – Hake, November 2014*
 734 *(EVHOE; commercial data from otter bottom trawls targeting demersal species*
 735 *OTB_DEF). Middle – Sole, November 2017 (Orhago; commercial data from otter bottom*
 736 *trawls targeting demersal species OTB_DEF). Right – squid, year 2015 (EVHOE;*
 737 *commercial data from otter bottom trawls targeting cephalopods OTB_CEP).*

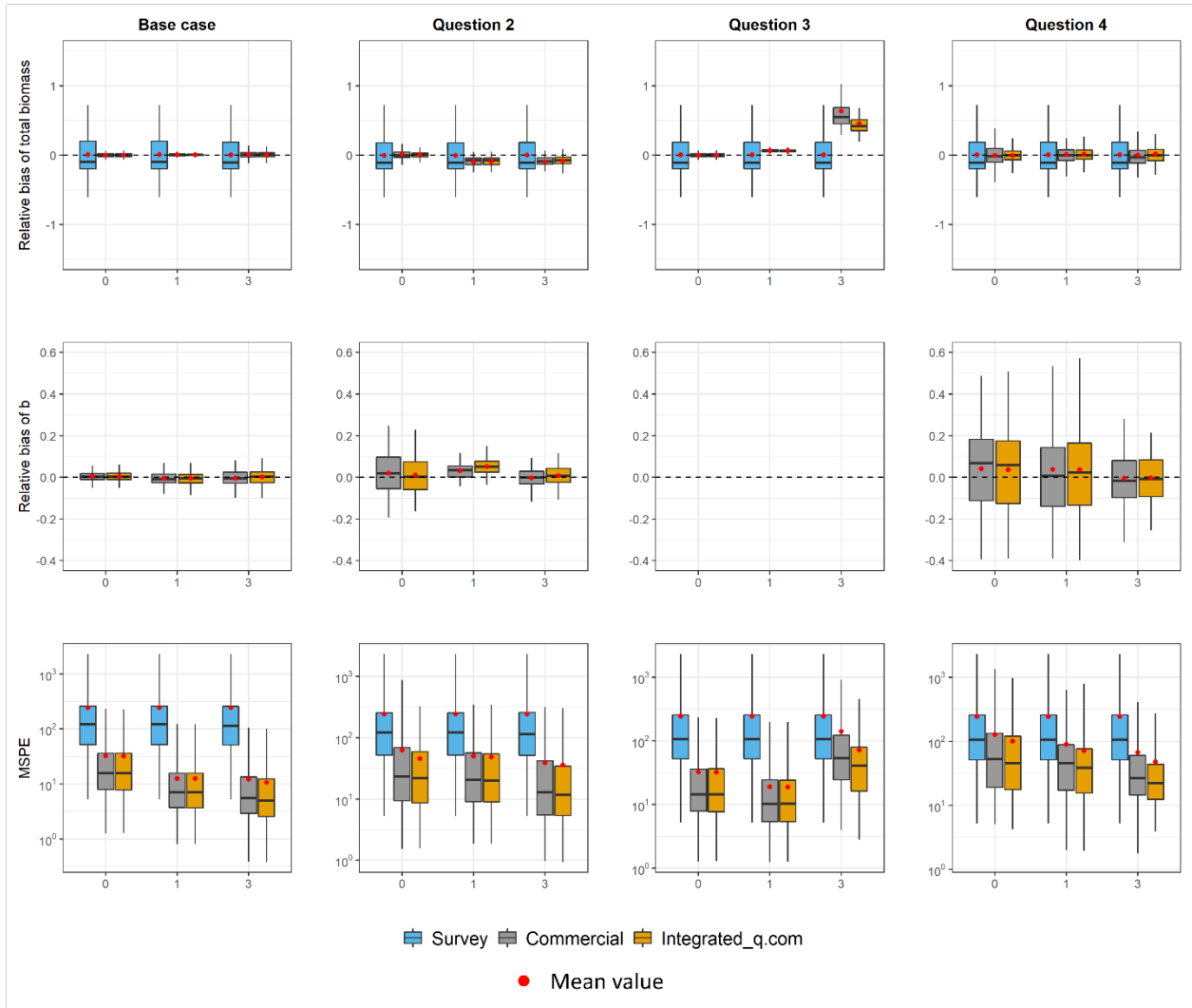


739
 740
 741
 742
 743
 744

Figure 4. Performance metrics obtained for various commercial and scientific data sample size. Column: intensity of the preferential sampling in simulated data. x-axis: 5 combinations of commercial and scientific sample size. 'com' stands for commercial, 'sci' stands for scientific, S stands for small sample size (50), M stands for middle sample size (400), L stands for large sample size (3000). Colors: model configurations.

745 *Integrated_q.com: integrated model with catchability fixed to 1 for commercial data;*
746 *Integrated_q.sci: integrated model with catchability fixed to 1 for scientific data. Boxplots*
747 *represent the variability among the 100 replicates.*

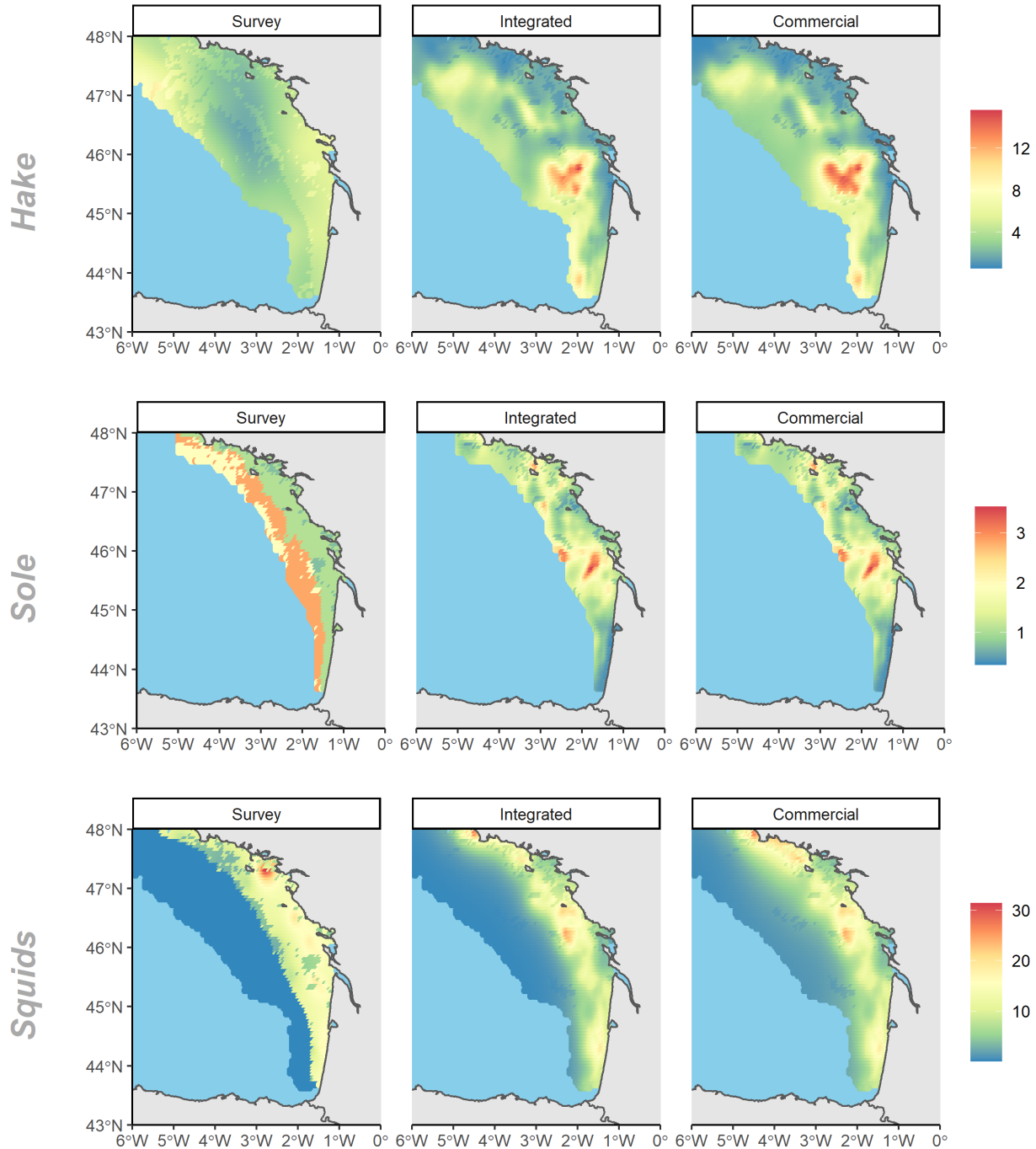
748



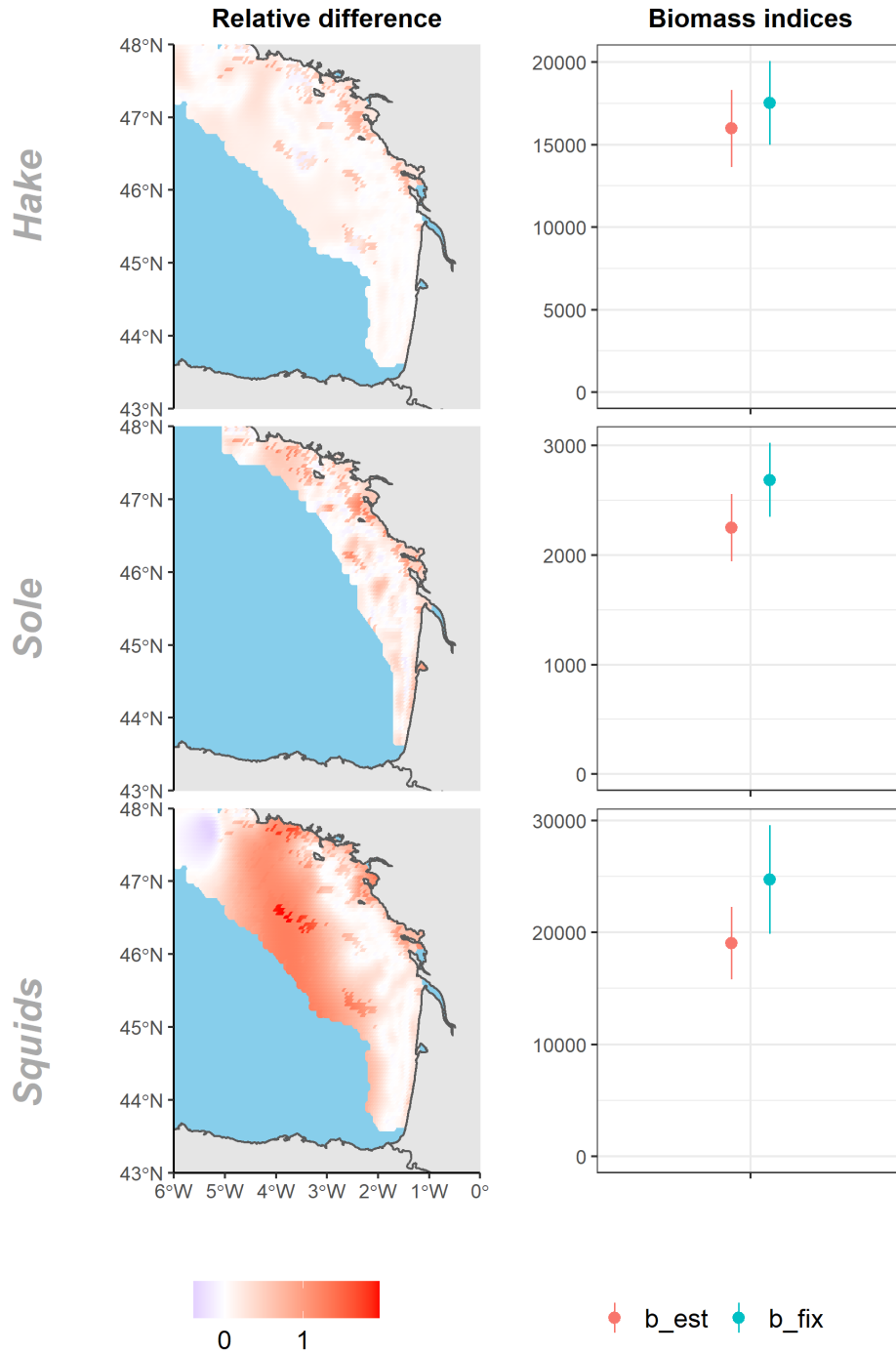
750

751

752 *Figure 5. Performance metrics obtained in different data and model configurations. Red*
 753 *points: mean value. 1st column: no discrepancy between simulation and estimation. 2nd*
 754 *column: commercial data do not cover a 9 x 9 zone of the grid. 3rd column: b is*
 755 *arbitrarily fixed to 0 in the estimation models. 4th column: data simulated with a random*
 756 *effect η in the sampling intensity process. In all configurations, simulations are*
 757 *conducted for three levels of preferential sampling (x-axis: $b = 0$, $b = 1$, $b = 3$). Colors:*
 758 *data sources used in the integrated model for inferences. Integrated_q.com: integrated*
 759 *model with catchability fixed with commercial data. Boxplots represent the variability*
 760 *among the 100 replicates.*
 761



763 *Figure 6. Prediction of the relative biomass for each case study. 1st column: model*
 764 *fitted to scientific data only; 2nd column: integrated model accounting for PS; 3rd*
 765 *column: commercial-based model accounting for PS. When the model is fitted to*
 766 *scientific data only, relative biomass is rescaled with the relative catchability parameter*
 767 *estimated within the integrated model so that all maps are in the same scale.*
 768



770 *Figure 7. Comparison of relative difference in biomass spatial predictions (calculated as*
 771 *$(S_{b_fix}(x) - S_{b_est}(x))/S_{b_est}(x)$ in space (left) and of total biomass (sum on the spatial*
 772 *domain; right) obtained with the integrated models from the 3 case studies when*
 773 *accounting or not for preferential sampling. b_est : PS is estimated. b_fix : PS is not*
 774 *accounted for.*
 775

Table 1: Simulations

General simulations description								
Biomass field	Depends on one continuous covariate (I_S) and one random spatial effect (δ). Both are simulated independently through a GRF with Matérn covariance function. Their range (ρ) and marginal variance are fixed respectively to 10 and 1. <small>n.b. the marginal variance quantifies the variability of the spatial process. For more details on marginal variance parameterization, see Lindgren <i>et al.</i> (2011).</small>						Simulated within a 25 x 25 grid.	
Scientific data	Random stratified plan within 4 strata (see Figure S2.1)			Catchability fixed to 1		Simulated with 10% of zeroes ($\xi_j = 0$)		
Commercial data	<p>Simulated according to three PS levels (i.e. three values for b - see Figure 2).</p> <ul style="list-style-type: none"> - $b = 0$: commercial sampling is not preferential; - $b = 1$: preferential sampling is moderate, commercial vessels mainly target areas where fish biomass is high; - $b = 3$: commercial sampling is highly preferential and vessels strongly target zones where biomass is high. <p>η is set to 0 for Q1, Q2, Q3. For Q4, η is set to tailor the sole case study. The range of η is set to 40 (4 times the range of δ), the marginal variance is set to 5 (5 times the marginal variance of δ). Catchability fixed to 1 Simulated with 30 % of zero when PS is null ($\xi_j = -1$)</p>							
	Simulation scenarios					Model configurations		
	b	Scientific sample size	Commercial samples size	Coverage of the study area	Additional random effect in sampling intensity (η)	Data sources considered in the model	PS estimated	Fixed catchability
Question 1: How do each data source contribute to inferences?	0,1,3	50	50, 400, 3000	Full	No	Scientific only, commercial only, both	yes	Scientific or Commercial
	0,1,3	50, 400, 3000	3000	Full	No	Scientific only, commercial only, both	yes	Scientific or Commercial
Question 2: How does a partial coverage of the study area by the commercial data affect the quality of the estimation?	0,1,3	50	3000	No fishing in a 9x9 cells box	No	Scientific only, commercial only, both	yes	Commercial
Question 3: What is the cost of ignoring PS in estimation when sampling is preferential?	0,1,3	50	3000	Full	No	Scientific only, commercial only, both	no (b fixed at 0)	Commercial
Question 4: How does the estimation perform when additional processes other than PS drive the fishing locations?	0,1,3	50	3000	Full	Yes	Scientific only, commercial only, both	yes	Commercial

779 **REFERENCES**

- 780
- 781 Abbott, J., Haynie, A., and Reimer, M. 2015. Hidden Flexibility: Institutions, Incentives, and
782 the Margins of Selectivity in Fishing. *Land Economics*, 91: 169–195.
- 783 Banerjee, S., Carlin, B. P., and Gelfand, A. E. 2014. Hierarchical modeling and analysis for
784 spatial data. CRC press.
- 785 Bourdaud, P., Travers-Trolet, M., Vermard, Y., and Marchal, P. 2019. Improving the
786 interpretation of fishing effort and pressures in mixed fisheries using spatial overlap
787 metrics. *Canadian Journal of Fisheries and Aquatic Sciences*, 76: 586–596.
- 788 Cameletti, M., Lindgren, F., Simpson, D., and Rue, H. 2013. Spatio-temporal modeling of
789 particulate matter concentration through the SPDE approach. *AStA Advances in
790 Statistical Analysis*, 97: 109–131.
- 791 Cheung, W. W., Lam, V. W., Sarmiento, J. L., Kearney, K., Watson, R., and Pauly, D. 2009.
792 Projecting global marine biodiversity impacts under climate change scenarios. *Fish
793 and fisheries*, 10: 235–251. Wiley Online Library.
- 794 Conn, P. B., Thorson, J. T., and Johnson, D. S. 2017. Confronting preferential sampling when
795 analysing population distributions: diagnosis and model-based triage. *Methods in
796 Ecology and Evolution*, 8: 1535–1546.
- 797 Cornou, A.-S., Quinio-Scavinner, M., Sagan, J., Cloâtre, T., Dubroca, L., and Billet, N. 2021.
798 Captures et rejets des métiers de pêche français - Résultats des observations à bord
799 des navires de pêche professionnelle en 2019. Ifremer.
- 800 Cressie, N. A. 1993. *Statistics for spatial data*. John Willy and Sons. Inc., New York.
- 801 Delage, N., and Le Pape, O. 2016. Inventaire des zones fonctionnelles pour les ressources
802 halieutiques dans les eaux sous souveraineté française. Première partie: définitions,
803 critères d'importance et méthode pour déterminer des zones d'importance à
804 protéger en priorité. Rapport de recherche. Pôle halieutique AGROCAMPUS OUEST,
805 Rennes.
- 806 Deporte, N., Ulrich, C., Mahévas, S., Demanèche, S., and Bastardie, F. 2012. Regional métier
807 definition: a comparative investigation of statistical methods using a workflow
808 applied to international otter trawl fisheries in the North Sea. *ICES Journal of Marine
809 Science*, 69: 331–342. Oxford University Press.
- 810 Diggle, P. J., Menezes, R., and Su, T. 2010. Geostatistical inference under preferential
811 sampling. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 59:
812 191–232.
- 813 Diggle, P. J. 2013. *Statistical analysis of spatial and spatio-temporal point patterns*. CRC
814 press.
- 815 Ducharme-Barth, N. D., Grüss, A., Vincent, M. T., Kiyofuji, H., Aoki, Y., Pilling, G., Hampton, J.,
816 *et al.* 2022. Impacts of fisheries-dependent spatial sampling patterns on catch-per-
817 unit-effort standardization: A simulation study and fishery application. *Fisheries
818 Research*, 246: 106169.
- 819 Erisman, B. E., Grüss, A., Mascareñas-Osorio, I., Licon-González, H., Johnson, A. F., and López-
820 Sagástegui, C. 2020. Balancing conservation and utilization in spawning aggregation

821 fisheries: a trade-off analysis of an overexploited marine fish. *ICES Journal of Marine*
822 *Science*, 77: 148–161. Oxford University Press.

823 Ferraris, J. 2002. Fishing fleet profiling methodology. Food & Agriculture Org.

824 Fletcher, R. J., Hefley, T. J., Robertson, E. P., Zuckerberg, B., McCleery, R. A., and Dorazio, R. M.
825 2019. A practical guide for combining data to model species distributions. *Ecology*,
826 100: e02710.

827 Francis, R. C. 2017. Revisiting data weighting in fisheries stock assessment models.
828 *Fisheries Research*, 192: 5–15.

829 Gelfand, A. E., Diggle, P., Guttorp, P., and Fuentes, M. 2010. Handbook of spatial statistics.
830 CRC press.

831 Gerritsen, H., and Lordan, C. 2010. Integrating vessel monitoring systems (VMS) data with
832 daily catch data from logbooks to explore the spatial distribution of catch and effort
833 at high resolution. *ICES Journal of Marine Science*, 68: 245–252.

834 Gimenez, O., Buckland, S. T., Morgan, B. J. T., Bez, N., Bertrand, S., Choquet, R., Dray, S., *et al.*
835 2014. Statistical ecology comes of age. *Biology Letters*, 10: 20140698. Royal Society.

836 Girardin, R., Hamon, K. G., Pinnegar, J., Poos, J. J., Thébaud, O., Tidd, A., Vermard, Y., *et al.*
837 2017. Thirty years of fleet dynamics modelling using discrete-choice models: What
838 have we learned? *Fish and Fisheries*, 18: 638–655.

839 Grüss, A., Thorson, J. T., Carroll, G., Ng, E. L., Holsman, K. K., Aydin, K., Kotwicki, S., *et al.*
840 2020. Spatio-temporal analyses of marine predator diets from data-rich and data-
841 limited systems. *Fish and Fisheries*, 21: 718–739.

842 Guisan, A., and Zimmermann, N. E. 2000. Predictive habitat distribution models in ecology.
843 *Ecological modelling*, 135: 147–186. Elsevier.

844 Haynie, A. C., Hicks, R. L., and Schnier, K. E. 2009. Common property, information, and
845 cooperation: Commercial fishing in the Bering Sea. *Ecological Economics*, 69: 406–
846 413.

847 Hilborn, R., and Walters, C. J. (Eds). 1992. Quantitative Fisheries Stock Assessment: Choice,
848 Dynamics and Uncertainty. Springer US.
849 <https://www.springer.com/gp/book/9780412022715> (Accessed 14 June 2021).

850 Hilborn, R., and Walters, C. J. 2013. Quantitative Fisheries Stock Assessment: Choice,
851 Dynamics and Uncertainty. Springer Science & Business Media. 575 pp.

852 Hintzen, N. T., Bastardie, F., Beare, D., Piet, G. J., Ulrich, C., Deporte, N., Egekvist, J., *et al.* 2012.
853 VMStools: Open-source software for the processing, analysis and visualisation of
854 fisheries logbook and VMS data. *Fisheries Research*, 115: 31–43. Elsevier.

855 Hintzen, N. T. 2021. Zooming into small-scale fishing patterns: The use of vessel monitoring
856 by satellite in fisheries science. Wageningen University.

857 Hintzen, N. T., Aarts, G., Poos, J. J., Van der Reijden, K. J., and Rijnsdorp, A. D. 2021.
858 Quantifying habitat preference of bottom trawling gear. *ICES Journal of Marine*
859 *Science*, 78: 172–184.

860 ICES. 2005. Report of the Workshop on Survey Design and Data Analysis (WKSAD). Sète,
861 France.

862 ICES. 2012. Manual for the international bottom trawl surveys. SISP 1-IBTS Copenhagen,
863 Denmark.

864 ICES. 2020a. International Bottom Trawl Survey Working Group (IBTSWG). ICES Scientific
865 Reports. ICES. <http://www.ices.dk/sites/pub/Publication>
866 Reports/Forms/DispForm.aspx?ID=37066 (Accessed 28 May 2021).

867 ICES. 2020b. Working Group for the Bay of Biscay and the Iberian Waters Ecoregion
868 (WGBIE). ICES Scientific Reports. ICES. <http://www.ices.dk/sites/pub/Publication>
869 [Reports/Forms/DispForm.aspx?ID=36841](http://www.ices.dk/sites/pub/Publication) (Accessed 28 May 2021).

870 Kai, M., Thorson, J. T., Piner, K. R., and Maunder, M. N. 2017. Spatiotemporal variation in
871 size-structured populations using fishery data: an application to shortfin mako
872 (*Isurus oxyrinchus*) in the Pacific Ocean. *Canadian Journal of Fisheries and Aquatic*
873 *Sciences*, 74: 1765–1780.

874 Kristensen, K., Thygesen, U. H., Andersen, K. H., and Beyer, J. E. 2014. Estimating spatio-
875 temporal dynamics of size-structured populations. *Canadian Journal of Fisheries and*
876 *Aquatic Sciences*, 71: 326–336.

877 Kristensen, K., Nielsen, A., Berg, C. W., Skaug, H., and Bell, B. M. 2016. TMB: Automatic
878 Differentiation and Laplace Approximation. *Journal of Statistical Software*, 70: 1–21.

879 Le Pape, O., Chauvet, F., Mahévas, S., Lazure, P., Guérault, D., and Désaunay, Y. 2003.
880 Quantitative description of habitat suitability for the juvenile common sole (*Solea*
881 *solea*, L.) in the Bay of Biscay (France) and the contribution of different habitats to
882 the adult population. *Journal of Sea Research*, 50: 139–149.

883 Lindgren, F., Rue, H., and Lindström, J. 2011. An explicit link between Gaussian fields and
884 Gaussian Markov random fields: the stochastic partial differential equation
885 approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*,
886 73: 423–498. Wiley Online Library.

887 Martínez-Minaya, J., Cameletti, M., Conesa, D., and Pennino, M. G. 2018. Species distribution
888 modeling: a statistical review with focus in spatio-temporal issues. *Stochastic*
889 *environmental research and risk assessment*, 32: 3227–3244. Springer.

890 Maureaud, A., Frelat, R., Pécuchet, L., Shackell, N., Mérigot, B., Pinsky, M. L., Amador, K., *et al.*
891 2020. Are we ready to track climate-driven shifts in marine species across
892 international boundaries? - A global survey of scientific bottom trawl data. *Global*
893 *Change Biology: gcb.15404*.

894 Moriarty, M., Sethi, S. A., Pedreschi, D., Smeltz, T. S., McGonigle, C., Harris, B. P., Wolf, N., *et al.*
895 2020. Combining fisheries surveys to inform marine species distribution modelling.
896 *ICES Journal of Marine Science*, 77: 539–552. Oxford University Press.

897 Murray, L. G., Hinz, H., Hold, N., and Kaiser, M. J. 2013. The effectiveness of using CPUE data
898 derived from Vessel Monitoring Systems and fisheries logbooks to estimate scallop
899 biomass. *ICES Journal of Marine Science*, 70: 1330–1340.

900 Nielsen, J. R. 2015. Methods for integrated use of fisheries research survey information in
901 understanding marine fish population ecology and better management advice:
902 improving methods for evaluation of research survey information under
903 consideration of survey fish detection and catch efficiency. Wageningen University.

904 Ocean Studies Board, and National Research Council. 2000. Improving the collection,
905 management, and use of marine fisheries data. National Academies Press.

906 Okamura, H., Morita, S. H., Funamoto, T., Ichinokawa, M., and Eguchi, S. 2018. Target-based
907 catch-per-unit-effort standardization in multispecies fisheries. *Canadian Journal of*
908 *Fisheries and Aquatic Sciences*, 75: 452–463.

909 Paradinas, I., Conesa, D., Pennino, M., Muñoz, F., Fernández, A., López-Quílez, A., and Bellido,
910 J. 2015. Bayesian spatio-temporal approach to identifying fish nurseries by
911 validating persistence areas. *Marine Ecology Progress Series*, 528: 245–255.

912 Parent, E., and Rivot, E. 2012. Introduction to hierarchical Bayesian modeling for ecological
913 data. CRC Press.

914 Pati, D., Reich, B. J., and Dunson, D. B. 2011. Bayesian geostatistical modelling with
915 informative sampling locations. *Biometrika*, 98: 35–48.

916 Pelletier, D., and Ferraris, J. 2000. A multivariate approach for defining fishing tactics from
917 commercial catch and effort data. *Canadian Journal of Fisheries and Aquatic
918 Sciences*, 57: 51–65. NRC Research Press.

919 Pennino, M. G., Conesa, D., Lopez-Quilez, A., Munoz, F., Fernández, A., and Bellido, J. M. 2016.
920 Fishery-dependent and-independent data lead to consistent estimations of essential
921 habitats. *ICES Journal of Marine Science*, 73: 2302–2310. Oxford University Press.

922 Pennino, M. G., Paradinas, I., Illian, J. B., Muñoz, F., Bellido, J. M., López-Quílez, A., and Conesa,
923 D. 2019. Accounting for preferential sampling in species distribution models.
924 *Ecology and evolution*, 9: 653–663.

925 Peterson, C. D., Courtney, D. L., Cortés, E., and Latour, R. J. 2021. Reconciling conflicting
926 survey indices of abundance prior to stock assessment. *ICES Journal of Marine
927 Science*, 78: 3101–3120.

928 Planque, B., Loots, C., Petitgas, P., LINDSTRØM, U. L. F., and Vaz, S. 2011. Understanding what
929 controls the spatial distribution of fish populations using a multi-model approach.
930 *Fisheries Oceanography*, 20: 1–17.

931 Punt, A. E. 2017. Some insights into data weighting in integrated stock assessments.
932 *Fisheries Research*, 192: 52–65.

933 Punt, A. E., Dunn, A., Elvarsson, B. Þ., Hampton, J., Hoyle, S. D., Maunder, M. N., Methot, R. D.,
934 *et al.* 2020. Essential features of the next-generation integrated fisheries stock
935 assessment package: A perspective. *Fisheries Research*, 229: 105617.

936 R Core Team. 2020. R: A language and environment for statistical computing. R Foundation
937 for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>.

938 Rivoirard, J., Simmonds, J., Foote, K. G., Fernandes, P., and Bez, N. 2008. Geostatistics for
939 estimating fish abundance. John Wiley & Sons.

940 Rochette, S., Rivot, E., Morin, J., Mackinson, S., Riou, P., and Le Pape, O. 2010. Effect of
941 nursery habitat degradation on flatfish population: Application to *Solea solea* in the
942 Eastern Channel (Western Europe). *Journal of sea Research*, 64: 34–44. Elsevier.

943 Rufener, M.-C., Kristensen, K., Nielsen, J. R., and Bastardie, F. 2021. Bridging the gap between
944 commercial fisheries and survey data to model the spatiotemporal dynamics of
945 marine species. *Ecological Applications*: e02453.

946 Salas, S., and Gaertner, D. 2004. The behavioural dynamics of fishers: management
947 implications. *Fish and Fisheries*, 5: 153–167.

948 Saunders, S. P., Farr, M. T., Wright, A. D., Bahlai, C. A., Ribeiro Jr., J. W., Rossman, S., Sussman,
949 A. L., *et al.* 2019. Disentangling data discrepancies with integrated population
950 models. *Ecology*, 100: e02714.

951 Schaub, M., and Abadi, F. 2011. Integrated population models: a novel analysis framework
952 for deeper insights into population dynamics. *Journal of Ornithology*, 152: 227–237.
953 Springer.

954 Schmitten, R. A. 1999. Essential fish habitat: opportunities and challenges for the next
955 millennium. *In American Fisheries Society Symposium*, p. 10.

956 Stephens, A., and MacCall, A. 2004. A multispecies approach to subsetting logbook data for
957 purposes of estimating CPUE. *Fisheries Research*, 70: 299–310.

- 958 Thorson, J. T., and Ward, E. J. 2014. Accounting for vessel effects when standardizing catch
959 rates from cooperative surveys. *Fisheries Research*, 155: 168–176.
- 960 Thorson, J. T., Ianelli, J. N., Munch, S. B., Ono, K., and Spencer, P. D. 2015a. Spatial delay-
961 difference models for estimating spatiotemporal variation in juvenile production and
962 population abundance. *Canadian journal of fisheries and aquatic sciences*, 72: 1897–
963 1915.
- 964 Thorson, J. T., Scheuerell, M. D., Shelton, A. O., See, K. E., Skaug, H. J., and Kristensen, K.
965 2015b. Spatial factor analysis: a new tool for estimating joint species distributions
966 and correlations in species range. *Methods in Ecology and Evolution*, 6: 627–637.
967 Wiley Online Library.
- 968 Thorson, J. T. 2015. Spatio-temporal variation in fish condition is not consistently explained
969 by density, temperature, or season for California Current groundfishes. *Marine
970 Ecology Progress Series*, 526: 101–112.
- 971 Thorson, J. T., Fonner, R., Haltuch, M. A., Ono, K., and Winker, H. 2016. Accounting for
972 spatiotemporal variation and fisher targeting when estimating abundance from
973 multispecies fishery data. *Canadian Journal of Fisheries and Aquatic Sciences*, 74:
974 1794–1807.
- 975 Thorson, J. T., Jannot, J., and Somers, K. 2017a. Using spatio-temporal models of population
976 growth and movement to monitor overlap between human impacts and fish
977 populations. *Journal of Applied Ecology*, 54: 577–587.
- 978 Thorson, J. T., Johnson, K. F., Methot, R. D., and Taylor, I. G. 2017b. Model-based estimates of
979 effective sample size in stock assessment models using the Dirichlet-multinomial
980 distribution. *Fisheries Research*, 192: 84–93.
- 981 Thorson, J. T. 2018. Three problems with the conventional delta-model for biomass
982 sampling data, and a computationally efficient alternative. *Canadian Journal of
983 Fisheries and Aquatic Sciences*, 75: 1369–1382. NRC Research Press.
- 984 Thorson, J. T., Adams, G., and Holsman, K. 2019. Spatio-temporal models of intermediate
985 complexity for ecosystem assessments: A new tool for spatial fisheries management.
986 *Fish and Fisheries*, 20: 1083–1099.
- 987 Thorson, J. T., Cunningham, C. J., Jorgensen, E., Havron, A., Hulson, P.-J. F., Monnahan, C. C.,
988 and von Szalay, P. 2021. The surprising sensitivity of index scale to delta-model
989 assumptions: Recommendations for model-based index standardization. *Fisheries
990 Research*, 233: 105745.
- 991 Trenkel, V. M., Beecham, J. A., Blanchard, J. L., Edwards, C. T., and Lorange, P. 2013. Testing
992 CPUE-derived spatial occupancy as an indicator for stock abundance: application to
993 deep-sea stocks. *Aquatic living resources*, 26: 319–332. EDP Sciences.
- 994 Winker, H., Kerwath, S. E., and Attwood, C. G. 2013. Comparison of two approaches to
995 standardize catch-per-unit-effort for targeting behaviour in a multispecies hand-line
996 fishery. *Fisheries Research*, 139: 118–131.
- 997 Witman, J. D., and Roy, K. 2009. *Marine Macroecology*. University of Chicago Press. 442 pp.
- 998 Zipkin, E. F., Inouye, B. D., and Beissinger, S. R. 2019. Innovations in data integration for
999 modeling populations. *Ecology*, 100: e02713.

1000