

1 **Why do lake whitefish move long distances in Lake Huron? Bayesian variable**
2 **selection of factors explaining fish movement distance**

3
4 Yang Li ^{a*}, James R. Bence ^a, Zhen Zhang ^b, and Mark P. Ebener ^c

5
6 ^a Quantitative Fisheries Center, Department of Fisheries and Wildlife, Michigan State University, East
7 Lansing, MI 48824, USA

8 ^b Enabling Capabilities Technology Center, Dow AgroSciences, Indianapolis, IN 46268, USA

9 ^c Inter-Tribal Fisheries and Assessment Program, Chippewa Ottawa Resource Authority, Sault Ste. Marie,
10 MI 49783, USA

11
12 **Abstract**

13 Understanding fish movement patterns is vital for stock assessment and fishery
14 management. We used a variable selection procedure in a Bayesian framework to
15 understand what factors most likely affect the net movement distance of individual fish
16 based on a conventional tag-recovery study of lake whitefish populations in Lake Huron
17 during 2003-2011, where fish of this species with spawning site fidelity were tagged
18 during the spawning season and recovered throughout the year. We found that fish with
19 greater total length, and those that were tagged and released from tagging sites near
20 Cheboygan and Alpena, Michigan, moved longer net distances than fish from other
21 tagging sites. Habitat conditions also had a profound effect on net movement distance.
22 We found that shorter movement distances by lake whitefish can be expected if the
23 relative density of the benthic amphipod *Diporeia* spp. was higher near the tagging site
24 during the recovery year. We also found evidence that lake whitefish may start their
25 annual spawning migration runs earlier during warmer years. More generally, our
26 Bayesian framework for analysis of conventional tagging data has potential for wide
27 applicability, and model details and our code are provided to facilitate this.

28
29 *Corresponding author: E-mail address: liyong11@msu.edu (Y. Li); Tel +1 (517) 355-
30 0126.

31
32 Key words: Fish Movement, Variable Selection, Tag Recovery, Lake Whitefish,
33 Bayesian.

34

35 **1. Introduction**

36 Many fish species move for long distances at various times during their life cycle, and
37 movements made by individuals vary from regular and predictable migration to less-
38 predictable resource driven nomadism (Runge et al., 2014). Most previous research that
39 evaluated changes in fish spatial locations focused on either the triggering factors or
40 distance between initial and final fish location (e.g., Albanese et al., 2004; Radinger and
41 Wolter, 2014), or on estimating net movement/migration rates of populations (Polacheck
42 et al., 2006; Vandergoot and Brenden, 2014).

43 Fish movement is essential from both conservation and management perspectives.
44 Movement behavior can influence how fish are distributed, whether their populations
45 persist in the face of ecosystem changes, and how stocks are assessed. Fish movement
46 can further influence ecological interactions and evolution (Lidicker and Stenseth, 1992).
47 Management problems such as inaccurate assessment results, or inappropriate catch
48 limits, can occur when actual fish movements do not agree with the spatial assumptions
49 made in stock assessments and management decisions, which can result in local
50 population depletion and population collapse (Fu and Fanning, 2004; Hutchings, 1996; Li
51 et al., 2015; Mitchell and Beauchamp, 1988; Rothschild, 2007).

52 Despite its ecological and management importance, understanding of fish movement
53 patterns in time and space, and how movements are related to environmental variables, is
54 still limited. Moreover, most previous research that focused on the triggering factors (i.e.,
55 factors causing the initiation of movement) and net fish movement distance were limited
56 to stream fish, given the easy calculation of net distance moved from conventional
57 tagging data. Much less is known about movement of fish that live in large water areas.
58 Most of which is known has been derived from electronic tagging data, although there are
59 many long-term conventional tagging programs. While technological advances make the
60 use of acoustic or pop-up tags increasingly useful, conventional tags are still more widely
61 used for estimating population size, mortality, and tracking individual growth, given their
62 lower price. Conventional tagging data can also provide information on the location at tag
63 release and tag recovery, which could be used for the estimation of movement route and
64 intensity (e.g., net fish movement distance) (e.g., Albanese et al., 2004; Gilliam and
65 Fraser, 2001).

66 The goal of this study was to develop a model framework for analysis of how factors
67 impact the distance fish move from when they are tagged until they are recovered ('net
68 fish movement distance' hereafter) in a larger water body, based on conventional tag-
69 recovery results. We based our research on several lake whitefish (*Coregonus*
70 *clupeaformis*) spawning stocks in Lake Huron of the Laurentian Great Lakes of North
71 America. As an ecological and economically important fish species in the Great Lakes,
72 lake whitefish have been found to move freely among multiple management units during
73 the non-spawning period, but show a high degree of natal homing, so nearly all mature
74 fish return to spawn at the same location each year (Ebener et al., 2010b). Previous
75 research on lake whitefish movement patterns provides a useful platform for us to derive
76 a priori hypotheses about the potential factors that influence movement. Since the
77 establishment of dreissenid mussels in the early 1990s, the ecosystem of four of the five
78 Great Lakes have changed substantially, including an overall decrease in the density of
79 lake whitefish's preferred food- *Diporeia* spp. (Barbiero et al., 2011; McNickle et al.,
80 2006; Mohr and Nalepa, 2005). In this context, Rennie et al. (2012) evaluated the

81 relationship between lake whitefish migration distance and growth rate, and found that
82 the least mobile population of lake whitefish was supported by a remnant *Diporeia* spp.
83 population. Ebener et al. (2010b) found that stock identity and season of recapture
84 affected net movement distance most strongly, while the influence of variables such as
85 sex, year, fish total length, and time at large was weaker. Although the role of
86 temperature has not been directly implicated in explaining patterns in the fish movement,
87 the association between lake whitefish harvest and surface water temperature suggested
88 that such a connection may exist (Price et al., 2003).

89 The pioneering studies of net movement distance used either a regression-tree based
90 approach or ANOVA models to test whether net movement distance varied significantly
91 in association with the factors they evaluated (e.g., Albanese et al., 2004; Ebener et al.,
92 2010b; Radinger and Wolter, 2014). Because some studies estimated the effects of
93 different factors as additive (i.e., causing a given distance change rather than a percentage
94 change in net movement distances), it is hard to generalize the results from studies with
95 different spatial and temporal scales. When jointly considering multiple factors and
96 continuous covariates, the ANOVA approach can provide only a rough picture of the
97 continuous relationship between net movement distance and explanatory factors. Thus, a
98 more thorough regression analysis is needed. The regression-tree based approach seeks to
99 approximate nonlinearity and interactions in the relationships between the net movement
100 distances and multiple factors by recursively partitioning the data points according to the
101 categorization of the factors (Ebener et al., 2010b). Such partitioning may have difficulty
102 in interpreting the effects, if the observations from the same tag or recovery area happen
103 to be separated into different branches of the tree. Some regression-tree applications have
104 partitioned data by site (i.e., different sites on different branches), and this can make it
105 difficult to develop a general understanding of movement (Ebener et al., 2010b). In
106 addition, although it is possible for regression-tree based approaches to rank or select
107 variables based on variable importance measures, they do not provide any further insight
108 of the uncertainty associated with their rankings or selections. Also information criteria,
109 such as Akaike's information criterion and the Bayesian information criterion, commonly
110 used as penalization terms for the number of parameters in model, are not applicable for
111 nonparametric tree-based models (Claeskens and Hjort, 2008).

112 We therefore considered a global linear regression model that accounts for joint
113 effects of multiple factors and the heterogeneity among sites, to study the relationship
114 between the net movement distance and individual factors. We further conducted a
115 variable selection procedure under a Bayesian framework to explore the plausibility of
116 alternative regression models that include various explanatory variables, and assess the
117 associated uncertainty. Bayesian variable selection treats the regression model itself as
118 random among all possible models with different sets of variables. Thus, it accounts for
119 model uncertainty in the overall assessment of uncertainty by making inferences on how
120 probable alternative models are after consideration of the data. The implementation of
121 Bayesian variable selection via the reversible jump Markov chain Monte Carlo
122 (rjMCMC) (Green, 1995) procedure is substantially more efficient in exploring the model
123 space than the traditional approaches such as all-subsets-regression (Woznicki et al.,
124 2016). While we believe our approach has substantial advantages over regression-tree
125 approaches, it could miss some nonlinear effects that could be identified by regression-

126 trees. Thus, as a check on robustness we compared our results with those from
127 regression-tree methods.

128 We considered how net distance moved from tagging to recapture locations changed
129 monthly and over years, and how this net movement pattern depended upon tagging
130 location. In addition, we considered how life history traits, namely total length, and sex,
131 and habitat features, namely *Diporeia* spp. density and water temperature, played a role
132 in these net movement patterns. Thus, the variables we considered as potential
133 explanatory factors in this study were tagging year, recovery year, recovery month,
134 year(s) between tag and recovery, fish total length, sex, tagging (spawning) site, and the
135 habitat variables based on *Diporeia* spp. density and growing degree days.

136 Our goal was to provide not only insight on how those factors influenced lake
137 whitefish movement in Lake Huron, but also a model framework for analyzing
138 movement mechanism based on conventional tagging data. Although Bayesian variable
139 selection in linear regression is a long-established approach (Mitchell and Beauchamp,
140 1988), it was rarely used in ecology or more specifically for uncovering explanations for
141 movements (Drouineau et al., 2017; Ethier et al., 2017). Drouineau et al. (2017) used a
142 Bayesian state-space model to analyze the effects of different environmental factors in
143 triggering migration of silver eel in fragmented rivers. Ethier et al. (2017) used Bayesian
144 models and variable selection to evaluate how environmental variables influenced
145 regional variation in population trends of Bobolink. Both studies used a mixture
146 distribution of priors (i.e., normal plus zero-inflation), which were estimated using a
147 Gibbs sampler. However, their variable selection procedure did not introduce a penalty
148 such as BIC for increasing number of selected variables. Also the Gibbs sampler usually
149 involves scanning all variables at each iteration, which could be computational
150 expensive, especially when the number of candidate variables is large.

151 To the best of our knowledge, this study is the first to apply the Bayesian variable
152 selection approach to compare the effects of various factors on fish net movement
153 distance by introducing an explicit prior penalty on model complexity, and the most
154 comprehensive to date in terms of the range of factors affecting whitefish movement. To
155 avoid sampling all indicators within a Gibbs sampler circle as in Drouineau et al. (2017)
156 and Ethier et al. (2017), we adopt the reversible jump MCMC algorithm for model
157 exploration that mimics stepwise selection and subsets regression technique, which is
158 more computationally efficient. Thus our research introduces an approach to fish
159 movement studies, which has the potential to be much more effectively interrogate a
160 large number of predictor variables. To facilitate usage of our approach, we provide the
161 open-source code for MATLAB program which is online available at to implement the
162 method.

163 **2. Methods**

164 *2.1 Data collection, selection, and calculation of net-movement distance*

165 Lake whitefish were tagged and released in a study coordinated by one of us (Mark P.
166 Ebener) at 21 individual tagging sites from nine spawning stocks in Lake Huron from late
167 October through December (i.e., spawning season) of 2003-2006. Total length (mm) of
168 all 35,285 tagged fish were measured before release, spatial coordinates of the tagging
169 and release location, and date of release were recorded for each fish. Lake whitefish
170 were tagged on or very near the spawning grounds and subsequently killed when
171 recovered by the commercial or recreational fishery. The commercial fishing season for

172 lake whitefish is not closed in Ontario waters during the spawning season, but it is closed
173 in Michigan waters. Thus, fish tagged and released at Detour, Cheboygan, Alpena, and
174 Saginaw Bay (Fig. 1) were extant 1-4 weeks before being subjected to fishing and tag
175 recovery. At Burnt Island, the Fishing Islands, and Sarnia fish were also tagged during
176 the spawning season, but commercial fishing was occurring simultaneously during
177 tagging so they had little time to be extant prior to tag recovery. Recovery happened from
178 December 2003 until December 2012, with the majority being recovered by commercial
179 fishermen, and the rest recovered during fishery surveys. Subsets of the data used here
180 were previously reported by Ebener et al. (2010a, 2010b), and details of the tagging
181 methodology are given by Ebener et al. (2010a).

182 Our analysis focused on drivers of net movement distance of lake whitefish tagged
183 and recovered in Lake Huron. We thus restricted attention to recoveries for which net
184 distance movement could be calculated and for which explanatory variable data were
185 available. Only recoveries that had location information recorded (either by latitude and
186 longitude or by 10-minute by 10-minute statistical grid, treated as though recovered at the
187 grid center) were considered. In addition, we excluded observations from fish that were
188 recovered within two days of release, as well as those without their recapture date, sex, or
189 total length recorded (i.e., explanatory variables). We also removed fish that were
190 recovered from Lake Michigan because of our focus on movement within Lake Huron
191 and because our explanatory variables were from Lake Huron. We further excluded
192 recoveries from two tagging sites that each produced only two total recoveries, and the
193 two fish recovered during 2012. Thus of the total of 2,098 reported lake whitefish
194 recoveries, 1,368 recoveries were used in this study. Details of data exclusion are
195 described in Supplementary Table S1. These recovered fish had total lengths between
196 375-667 mm at the time of tagging, and were tagged and released from seven spawning
197 sites (Fig. 1).

198 We used log-transformed net movement distance as a response variable because net
199 movement distances were highly skewed. We calculated net movement distance based on
200 the shortest water distance between tagging and recovery locations, using a Dijkstra type
201 shortest path algorithm (Vincenty, 1975; online Appendix A). We standardized log-
202 transformed net movement distance by subtracting the mean and dividing by standard
203 deviation prior to analysis.

204 2.2 Explanatory variables

205 We hypothesized that net movement distance for lake whitefish in Lake Huron would
206 be influenced by 1) life history traits, which included total length, and sex; 2) temporal
207 factors, which included tagging year (tag_Y), recovery year (rec_Y), recovery month
208 (rec_M), and year(s) between tagging and recovery (year_lag); and 3) habitat condition,
209 which included *Diporeia* spp. density, and growing degree days; and 4) tagging
210 (spawning) sites. These hypotheses, related variables, and the expected sign of the
211 associated coefficients, if hypotheses were supported, are in Table 1. Due to the strong
212 spawning site fidelity of lake whitefish (i.e., nearly all lake whitefish move back to where
213 they born each year during the spawning season), we only considered the habitat
214 conditions during the recovery year as a predictor. That is, the net movement is in
215 actuality the net movement since the prior spawning season. We used relative *Diporeia*
216 spp. density, which was the *Diporeia* spp. density of the release location divided by the
217 mean of all sampled stations in Lake Huron for that year. The U.S. EPA Great Lakes

218 National Program Office collected *Diporeia* samples every August since 1999 at 12 Lake
219 Huron stations (Barbiero et al., 2011). The release location density was defined as the
220 density at the sampled location closest to the release location. Our hypothesis was that
221 lake whitefish tended to stay near their tagging locations when *Diporeia* density was
222 higher in that vicinity.

223 We proposed two alternative hypotheses for the relationship between growing degree
224 days (GDDs) (i.e., also known as thermal time, a weather-based indicator about heat
225 acumination for assessing fish growth; e.g., Chezik et al., 2014) and lake whitefish net
226 movement distance, and these led to two distinct sets of GDD variables. These two sets
227 were used in two alternative analyses. We calculated GDDs based on mean daily
228 (daytime) surface temperatures from the Great Lakes Surface Environmental Analysis
229 (GLSEA) remote sensing surface water temperature data (See online Appendix A).

230 *Case 1 (GDD hypothesis 1)*—Lake whitefish respond to growing conditions they had
231 experienced during the current year. Thus, they would tend to be closer to their tagging
232 (spawning) site when the growing degree days (GDD) at the tagging location was greater
233 than the lake average GDD during that same time period. This led us to define the
234 explanatory variable relative GDD difference (“GDD_Diff”), calculated as: $GDD_Diff =$
235 $(GDD_{tag} - GDD_{lake})/GDD_{lake}$, where GDD_{tag} and GDD_{lake} are the cumulated non
236 negative degree days (°C. days) that exceeded 5°C (Rennie et al., 2009) at the tagging
237 location or for the lake-average, respectively, from the first day of the recovery year to
238 the day of recovery.

239 *Case 2 (GDD hypothesis 2)*—The spawning season of lake whitefish would be shifted
240 earlier in the year, in years for which GDDs accumulated faster, because individual fish
241 would reach a physiological status allowing spawning earlier under such conditions.
242 Preliminary model fits without a GDD effect indicated that lake whitefish were generally
243 closer to the spawning location during September through December, than at other times
244 of the year. We therefore assumed that GDD might potentially influence net movement
245 distance (to varying degrees) only during these months. Thus, we added four additional
246 interaction variables (recovery month \times GDD_{lake}) for September through December
247 recoveries. We used GDD_{lake} because fish would be living and feeding away from their
248 spawning/tagging sites until moving to those sites for spawning.

249 After creating dummy variables and choosing the category with the largest number of
250 observations as the baseline category for each factor, we have a total of 34 (for GDD
251 hypothesis 1 case) and 37 (for GDD hypothesis 2 case) candidate variables including the
252 intercept (Table 1). Note that there was no dummy variable created for the baseline
253 category (i.e., tagging site: Detour, recovery month: June, tagging year: 2004, recovery
254 year: 2006, or sex: Male), because it was defined as zero for all other categories for that
255 factor. All explanatory variables were standardized like net movement distance.

256 2.3 Model framework

257 We used Bayesian variable selection to identify the highly probable subsets of
258 predictors for the linear regression and, given a set of predictors, we assessed likely
259 parameter values. Given the Bayesian approach we used, inferences were based on a
260 posterior distribution, which depends jointly on assumed prior distributions and the
261 likelihood of the data. Model components (i.e., regression model, prior distributions, and
262 likelihood) are described in Section 2.31 (Model Description) and how we used Markov

263 chain Monte Carlo (MCMC) techniques to derive posterior distributions in Section 2.3.2.
264 A separate model selection process was conducted for the two cases (GDD hypotheses).

265 2.3.1 Model description

266 Each possible model is of the form:

$$267 \quad Y = X_j \beta_j + \epsilon, \epsilon \sim N(0, \tau^2 I_N) \quad (\text{Eq.1})$$

268 where Y is the response variable (i.e., log-transformed net movement distance) with N
269 observations, X_j is the $N \times q_j$ design matrix (containing data for the predictors included
270 for that regression), β_j is a vector of parameter coefficients (an intercept included in
271 every model plus q_j-1 additional coefficients for the predictor variables included for that
272 model) and ϵ is the residual error. We assumed here homogenous, normal and
273 independent residual errors, with variance τ^2 . We assumed independent errors given the
274 relatively large distances between tagging sites (Fig. 1) and because tagging and recovery
275 spatial factors were included as potential explanatory variables. As described in online
276 Appendix B, I_N could be replaced with a selected correlation matrix. A model with a
277 specific subset of selected variables is represented by \mathcal{J} , which formally is an index set,
278 that maps the q variables in the selected model to the larger set of p possible variables.

279 The $\beta_j = (\beta_1, \beta_2, \dots, \beta_q)'$ had a normal prior $\beta_j \sim N(0, \tau^2 \Lambda)$, where $\Lambda =$
280 $\text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_q\}$. The λ_j s were modeled as arising from a higher level inverse-gamma
281 prior distribution ('hyperprior') with shape parameter a_λ and scale parameter b_λ . We
282 assumed a hyperprior with inverse-gamma density for τ^2 with shape a_τ and scale b_τ . The
283 hyperparameters were set to the values $a_\lambda = a_\tau = 2$ and $b_\lambda = b_\tau = 0.001$, which
284 correspond to a rather dispersed prior distribution.

285 The normal prior with diagonal variance-covariance matrix for the λ_j s represents a
286 decision to use a Bayesian counterpart to Ridge regression. The λ_j s represent the signal
287 to noise ratio of the effects in the model, and their magnitude played a role in whether an
288 effect was included and the size of selected models. Modeling them as arising from a
289 hyperprior (rather than specifying their values) allowed for adaptive learning on which
290 variables to include during the model search process.

291 We included an intercept in all models to account for the grand mean level of Y , as is
292 often done for variable selection. There are a total of 2^{p-1} possible models (i.e., an
293 intercept-only model, all possible models with one additional variable, all possible
294 models with two additional variables, etc.). We specified the prior probability of each
295 model as arising from the product of a prior probability for a model of a given size (i.e.,
296 $\pi(q)$), multiplied by the probability of a specific model given its size:

$$297 \quad \pi(\mathcal{J}) = \pi(\mathcal{J}_q, q) = \pi(\mathcal{J}_q | q) \pi(q) \quad (\text{Eq.2})$$

298 We let $\pi(q) \propto \exp\{-\kappa q\}$ for integer q from $\{1, 2, \dots, p\}$. Here \propto means
299 "proportional to" up to a constant that is irrelevant in making inferences about the hyper-
300 parameter κ . This placed higher prior probability on models with smaller size, as is
301 consistent with common practice in variable selection, and the rate at which the prior
302 probability falls as model size increases was determined by κ . We set $\kappa = \log(N) / 2$,
303 which is analogous to a BIC-type penalty on the number of selected variables (Schwarz,

304 1978). Conditional on q , each model J_q had an equal chance of being selected, i.e.,
305 $\pi(J_q|q) = 1/\binom{p-1}{q-1}$ for $q > 1$, and for $q = 1$ no selection is needed.

306 2.3.2 Characterization of the posterior Distribution using MCMC

307 We used Markov chain Monte Carlo methods to determine the posterior distribution.
308 We used a hybrid reversible jump technique (rjMCMC), because it performs well when
309 selecting among different sets of variables, which involved trans-dimensional states of
310 Markov chain (Green, 1995; Woznicki et al., 2016). Our procedure involved running
311 multiple chains and combining converged portions of these into one set of "retained
312 samples." The retained samples were summarized to highlight desired properties of the
313 posterior distributions.

314 Details on the implementation of the hybrid of rjMCMC for model search and Gibbs
315 sampler for parameters given the model, as well as procedures for evaluating MCMC
316 convergence and producing the retained samples are given in online Appendix B. We
317 summarized the posterior distributions for regression model parameters in two ways:

318 *Variable-wise summary*— This provided a summary conditional on the j -th variable
319 being selected. This was based on summarizing all samples included in the final MCMC
320 chains for a model that included the j^{th} variable. For the corresponding β_j , the posterior
321 mean and 95% (equal probability tail) credible intervals were constructed from these
322 samples. As a measure of the importance of each variable we also calculated the marginal
323 inclusion probability (Barbieri and Berger, 2004), as the proportion of all retained
324 MCMC samples that included the j^{th} variable in the model.

325 *Model-wise summary*— This was conditional on one specific model J in the posterior
326 samples, and thus was based only on retained MCMC samples for that model. We
327 provide such summaries for the 12 "top" models. Here models are ranked based on the
328 posterior probability, calculated as the proportion of all retained MCMC samples that
329 were model J . For the top models, we summarized the posterior distributions of the β_j s
330 for all variables in J , again in terms of the posterior mean and 95% credible interval.

331 2.3.3 Model diagnosis, simulation study, and comparison with tree-based methods

332 We used the posterior predictive assessment of model fitness using the χ^2 -
333 discrepancy (Gelman et al., 1996), based upon which we calculated the Bayesian p-value
334 for the top models in both GDD hypothesis cases. We also conducted simulations to
335 evaluate how well our Bayesian variable selection procedure can discover the true set of
336 important variables and estimate the corresponding effects, under five different scenarios
337 with varying combinations of true predictor variable effects. We also applied two tree-
338 based methods to our data, and compared the top variables from tree-based methods, the
339 gradient boosting regression tree method (Ethier et al., 2017) and the random forests
340 approach (Breiman, 2001), to the selected variables from our variable-wise summary.
341 Detailed methods for our diagnostic procedures, simulations, and tree-based applications
342 are given in online appendices C, D, and E respectively, and performance statistics
343 resulting from the simulations and tree-based methods are also presented in the
344 appendices.

345 **3. Results**

346 The posterior distributions of the number of selected variables were similar for the
347 two GDD cases and suggested that the most probable model sizes had 6 and 7 variables
348 including an intercept (Fig. 2). However, the selected variables were quite different (see
349 Section 3.1).

350 *3.1 Variable-wise summary*

351 *GDD hypothesis 1 Case:* There were 10 variables with 95% credible intervals that did
352 not cover 0, which we define as "consistent effects" (Fig. 3). Variables that had consistent
353 effects generally had high marginal inclusion probability, and more generally variables
354 with higher probability of inclusion tend to have more of their posterior distribution on
355 one side of zero (Fig. 3). The six top variables (length, tagging site: Cheboygan and
356 Alpena, *Diporeia*, and recovery months October and November) had marginal inclusion
357 probability above 0.75 (i.e., they are selected by more than 75% of the total posterior
358 samples). The variable recovery month September also had a relatively large marginal
359 inclusion probability (0.40). The other variables that were detected as consistent effects
360 had substantially lower marginal inclusion probability (<0.07) are: years lag, tagging site
361 Fishing islands, and recovery month December. According to the posterior mean of those
362 10 variables with consistent effects, fish with greater length, longer lag between the
363 tagging and recovery years, released at tagging site Cheboygan, Alpena, and Fishing
364 Islands, and recovered in December had greater net movement distance, while fish
365 released at the tagging site with higher density of *Diporeia*, and recovered during
366 September, October, and November had shorter net movement distance. Our first GDD
367 hypothesis was not supported by the variable selection results because the 95% credible
368 interval of the associated effect covered 0, and had a marginal inclusion probability of
369 only 0.004.

370 *GDD hypothesis 2 Case:* As was true for the previous case, consistency of effects and
371 the marginal probability of inclusion were positively associated (Fig. 4). The six top
372 variables in terms of marginal inclusion probability (length, tagging site: Cheboygan and
373 Alpena, *Diporeia*, recovery month November, interaction effect between lake-average
374 GDD and recovery month October) were similar to the top variables for GDD hypothesis
375 1. The major exceptions were for the recovery month October and the GDD associated
376 variables (Fig. 4). Consistent with the results for GDD hypothesis 1, fish with greater
377 length, longer lag between the tagging and recovery years, and released at tagging site
378 Cheboygan, Alpena, and Fishing islands had greater net movement distance, while fish
379 released at the tagging site with higher density of *Diporeia*, and recovered during
380 September, and November had shorter net movement distance. The effect of recovery
381 month October had a similar negative posterior mean, although the effect was less
382 consistent. The less consistent effect of October is likely associated with the inclusion of
383 the GDD associated variables for Hypothesis 2. Our second hypothesis of GDD was well
384 supported by the variable selection results. The interaction effect between lake-average
385 GDD and recovery month October, and the interaction effect between lake-average GDD
386 and recovery month November were both consistent, with a negative posterior mean and
387 the former was smaller than the latter. That is, fish tended to have smaller net movement
388 distance in October and November if the lake-average GDD was greater, and the effect
389 was larger in October than that in November. On the other hand, the interaction between
390 lake-average GDD and recovery month December was also consistent, but with a positive

391 posterior mean, which suggested that fish tended to have greater net movement distance
392 in December if the lake-average GDD was greater. An overall interpretation of these
393 effects is a shift in the spawning season in association with GDD, with more fish close to
394 the spawning grounds by October and having moved away by December, when GDD was
395 higher.

396 3.2 Model-wise summary

397 Given the variable selection result support our second GDD hypothesis, we only
398 present model-wise summary for GDD hypothesis 2 case. One or more of the four
399 interaction variables (recovery month \times GDD for the fish that were recovered from
400 September, October, November, and December) were included in at least one out of the
401 top 12 models. In addition to the three variables with marginal inclusion probability
402 equals 1 in Fig. 4 (*Diporeia*, tagging site: Cheboygan and Alpena), the interaction effect
403 recovery month October \times GDD was also included in all 12 top models (Fig. 5).
404 Recovery month November was included in nine of the 12 top models (all but models 5,
405 7, and 10), while the interaction variable November \times GDD_{lake} was included in the other
406 three top models. Total length of tagged fish was included in eight out of the 12 top
407 models, recovery month September was included in four out of the top models, the
408 interaction variable December \times GDD_{lake} was included in three out of the top models,
409 and recovery month December and the interaction variable September \times GDD_{lake} were
410 included in one out of the top models. The top two models both had a posterior
411 probability greater than 0.14. These models were similar. The best model (i.e., the highest
412 posterior probability model) included six variables and the second best model included all
413 those variables, plus recovery month September. Most estimated β_j s are consistent across
414 the top models, suggesting the effect of a variable was relatively uninfluenced by the
415 presence of other variables in the models.

416 The best model (Model 1 in Fig. 5) for the fit with GDD hypothesis 2 is summarized
417 in Table 1. From the best model, fish that were tagged and released from tagging sites
418 Cheboygan and Alpena had longer net distance than fish released at other tagging sites.
419 Lake whitefish with greater total length also tended to have greater net distance. Fish that
420 were recovered in November consistently had shorter net distance than fish recovered in
421 other months. In addition, shorter movement distance could be expected if the relative
422 *Diporeia* density was higher near the spawning locations during the recovery year. The
423 interaction term of month October and lake-average GDD resulted in shorter net distance
424 when lake-average GDD was high.

425 3.3 Model diagnosis, simulation study, and comparison with tree-based methods

426 *Model diagnosis*— There is no evidence for lack-of-fit of the top models under both
427 GDD hypothesis cases. In particular the scatterplot of predicted and realized χ^2 appear
428 consistent with a 1:1 relationship (Fig. S1 in online Appendix C) and the Bayesian p-
429 values are much larger than 0.05, indicating that the null hypothesis that the observed
430 data follow the hypothesized model is not rejected. We also did a residual analysis for the
431 top model of both GDD hypothesis cases, and plotted averaged standard residuals for the
432 MCMC samples associated with those top models versus selected (including both
433 continuous variables and two way combinations of categorical predictors). We did not
434 observe any suspicious patterns from the plot given: 1) all residuals are nearly symmetric

435 about zero, majority within (-3, 3), according to the 3-sigma rule, 2) there were no
436 obvious trends in variation or mean across different values of the predictors.

437 *Simulation study*— In general, our BVS method had consistent performance at
438 identifying important variables, and in identifying an appropriate model under scenarios
439 with varying combinations of candidate variables (see online Appendix D). Effects of
440 interactions, and of continuous and categorical variables were all likely to be selected
441 when they actually had effects, and not to be selected when they did not have effects on
442 the response variable. Across all scenarios, the true model was very likely to be included
443 in the top two models (i.e., probability ≥ 0.9), and most likely to be our top model (i.e.,
444 probability ≥ 0.74).

445 *Comparison with tree-based methods*— The top variables from both tree-based
446 methods in Figure S3 (online Appendix E) are consistent with Bayesian variable selection
447 (BVS) results, although there were several exceptions. The first exception was for the
448 GDD hypothesis 1 case, where GDD_Diff was not selected as important variable by
449 BVS, but was selected as top variables by both boosted regression tree and random forest
450 approaches. We believe that this is due to several high-leverage GDD_Diff observations
451 (Fig. S4 in online Appendix E), which the regression tree methods see as nonlinear
452 effects. A second exception, also for the GDD hypothesis 1 case, was that fish length had
453 a high inclusion probability (0.78 with BVS) and was also a top variable for boosted
454 regression trees but was not included in the top list for the random forests approach. A
455 third exception was that the rank of the variable September was lower for the tree-based
456 approaches than for BVS, and this was true for both GDD hypotheses, albeit the three
457 approaches rank variable importance in different ways (probability of inclusion for BVS,
458 see X axis of Fig. S3 and Fig. S5 for tree-based methods).

459 **4. Discussion**

460 The goal of this study was to develop a model framework for analysis of how factors
461 impact net fish movement distance in a larger water body, based on conventional tag-
462 recovery results, and apply the framework to lake whitefish spawning stocks in Lake
463 Huron of the Laurentian Great Lakes of North America. Our framework used a data-
464 driven Bayesian variable selection (BVS) method, where the candidate variables
465 represented hypothesis about drivers of net movement distance. The hypotheses we
466 evaluated were that the net movement distance of adult lake whitefish in the main basin
467 of Lake Huron was related to 1) fish total length, 2) sex, 3) tag and release year, 4)
468 recovery year, 5) recovery month, 6) year(s) between tagging and recovery, 7) *Diporeia*
469 spp. density near the spawning locations relative to the lake-wide *Diporeia* spp. density,
470 8) relative difference between the tagging site and lake-wide growing degree days, and 9)
471 the interaction term between lake-wide growing degree days and recovery month. Some
472 of the above hypotheses were well supported by the results presented.

473 There was a consistent positive relationship between lake whitefish net movement
474 distance and fish total length at the time of tagging. This is consistent with conclusions
475 from previous studies of stream-dwelling fish, in which longer movement and home
476 range was observed for larger fish (Gatz and Adams, 1994; Gunning and Shoop, 1963).
477 This greater movement may be due to the increasing mass-specific bioenergetic costs of
478 mobility with decreasing body size (Roff, 1991). Minns (1995) also found that the home
479 range is related to body size in freshwater fisheries and is consistently larger in lakes than
480 in rivers.

481 Because of the spawning site fidelity of lake whitefish, recovery months were
482 expected to have effects on net movement distance. Ebener et al. (2010b), analyzing
483 some of the same data but focused on different spatial and temporal scales with fewer
484 predictor variables, also demonstrated that season of recapture played an important role
485 in the distance moved by lake whitefish. Here, net movement distance was found to be
486 negatively related to recovery months September, October and November, and positively
487 related to December. This suggested that the spawning migration movement for lake
488 whitefish generally occurred within months from September to November, and after that,
489 fish tended to leave their spawning site and were actually further from the spawning
490 location than in the baseline month of June.

491 Past research has documented that some life history events such as reproduction can
492 be accelerated with warmer water temperature (Forseth et al., 1999). For example, the
493 spawning of walleye has occurred earlier with earlier ice-out related to warmer
494 temperature (Schneider et al., 2010). We found similar patterns in our study. When lake
495 average GDD was higher, lake whitefish tended to move or stay closer to their spawning
496 sites from September to November, and to be further away from their spawning sites in
497 December. This suggests that fish may start their annual spawning migration runs earlier
498 in warmer years after acquiring and processing energy needed for spawning. The
499 underlying mechanism could be that fish have to either achieve a critical condition before
500 the cost of migration/spawning can be offset (Forseth et al., 1999), or to accumulate
501 enough energy to survive a winter starvation period before spawning.

502 Although the decline of *Diporeia* spp. density in the Laurentian Great Lakes due to
503 the establishment of dreissenid mussels has been argued as the main reason of lake
504 whitefish expanding their movement range (Ebener et al., 2010b; Rennie et al., 2012), we
505 know of no other direct evaluation of an effect of *Diporeia* density on movement. Our
506 study evaluated this hypothesis by including relative *Diporeia* spp. density as a predictor
507 for lake whitefish net movement distance, and we found that when relative *Diporeia* spp.
508 density was high near the spawning grounds, lake whitefish tended to stay closer to their
509 spawning site. This implied that fish might expand their foraging area when *Diporeia*
510 density was low near their preferred habitat. Our analysis also found an effect of the
511 relative density of *Diporeia* within a year, which suggests a pattern related to the density
512 of this prey, not just a general change in movement over time throughout the lake as
513 *Diporeia* declined.

514 Lake whitefish tagged and released from the tagging sites Cheboygan and Alpena had
515 consistently greater net distance than those released from other areas. The underlying
516 reasons may be relate to the bathymetry and shoreline features of Lake Huron. Deep
517 water (>80 m) near Cheboygan and Alpena may restrict the movement of Cheboygan and
518 Alpena spawning stocks to north-south direction where there is a large area with relative
519 shallow water. In contrast, the spawning stocks in Detour and Burnt Inlands may be
520 constrained from moving south by the deep water in north of the main basin of Lake
521 Huron, so that they tended to move in the east-west direction. Considering the shape of
522 Lake Huron and the locations of those spawning stocks, movement in the north-south
523 direction allows longer movement distance than in the east-west direction.

524 There was similarity but also some differences in variable selection between our
525 Bayesian variable selection and tree-based methods. One notable difference between tree-
526 based methods and the Bayesian method is in the inclusion of GDD difference in GDD

527 hypothesis 1 case for the tree-based methods but not by BVS. The overall neutral effect
528 and low importance for the BVS was apparently because a few high-leverage points were
529 treated as noise. By recursively partitioning the data according to different ranges of
530 predictors, the tree-based methods are less sensitive to those points. However, such
531 localized results based on small samples can hardly provide any general predictability.
532 Rättsch et al. (2001) also found that overfitting can occur for regression tree-based
533 methods using a boosting algorithm when there is a lot of noise.

534 Our BVS method can be used for various different species and any water system
535 meeting our input requirements. For conventional tagging studies done in large lakes
536 (e.g., Lake Huron as in our case) or oceans, shortest water distance can be used as
537 response variable; while for a tagging study done in a river, a river network needs to be
538 built /considered for calculating (net) movement distance. Given that our Bayesian
539 variable selection method penalizes the number of selected variables, it has the potential
540 to perform well for other cases with more candidate explanatory variables than we used
541 in our application. In addition, the approach is adaptable to situations where residuals
542 might be correlated. We assumed no such correlations given the spatial distribution of
543 tagging sites and inclusion of spatial covariates (e.g., tagging sites), but in other situations
544 there could be spatial structure that should be accounted for in random part of the model.
545 In such cases correlations could be made a function of a measured quantity like distance
546 between tagging sites, and our code and detailed description of the model in the
547 supplement outlines how this can be done. In addition, our Bayesian method also allows
548 extra flexibility such as including: (1) random effects to cope with grouping variables
549 with a large number of outcomes, which can greatly improve the prediction by better
550 explaining the variability; (2) prior information for the effects of variables with flexible
551 choices that can be leveraged from a broad catalog in the Bayesian variable selection
552 literature. Thus we believe our work established a framework that could facilitate
553 additional studies of animal movement based on conventional tagging data.

554 We made several simplifying assumptions and choices in our analysis. Firstly, we
555 assumed 100% spawning site fidelity, so for the environmental factors *Diporeia* spp.
556 density and GDDs, only data for the year of recovery were used. While fidelity is likely
557 not 100%, available data suggest it is quite high for lake whitefish (Ebener et al., 2010b).
558 Secondly, the T_0 used for the calculation of cumulative GDD is 5°C (Rennie et al., 2009),
559 but it is possible that this is not the best threshold or that fish are responding to
560 temperature in a different or more complex fashion than we assumed. We believe that
561 violation of the 100% fidelity assumption and the GDD assumptions would act to obscure
562 effects of *Diporeia* and tagging site rather than cause us to discover artefactual effects.
563 Thirdly, we assumed similar tag reporting rates across all recovery basins, so data were
564 not weighted across different recovery basins. Violation of this assumption could be
565 influencing details of our results. However, we suspect the larger qualitative effects are
566 real rather than artifacts of such a violation. If there were dominating differences in tag
567 reporting rates among basin, we would have expected that to be reflected in consistent
568 tagging site effects for sites within basins, which we did not see in our results.

569 **Acknowledgements**

570 Funding for this project was provided by Quantitative Fisheries Center supporting
571 partners, Michigan DNR through the Partnership for Ecosystem Research and
572 Management program, and grants from the Great Lakes Fishery Commission Fishery

573 Research Program and the Great Lakes Fishery Trust (Project 2012.1250). We thank the
574 commercial fisherman and the staff of the Inter-Tribal Fisheries and Assessment
575 Program, the Hammond Bay Biological Station, the Great Lakes Fishery Commission,
576 and Michigan DNR for their helps in tagging fish in very uncomfortable weather
577 conditions. We thank Dr. Richard P. Barbiero for providing *Diporeia* abundances data,
578 and Dr. Lacey Mason for providing surface water temperature data. We acknowledge the
579 support of the Michigan State University High Performance Computing Center and the
580 Institute for Cyber-Enabled Research. We wish to thank the Modeling Subcommittee of
581 the Technical Fisheries Committee for 1836 Treaty-ceded waters for providing input on
582 our research. This is publication 2017-XX of the Michigan State University Quantitative
583 Fisheries Center.
584

585 **Appendix A. Calculation of shortest water distance and GDD**
586 Calculation of shortest water distance and GDD can be found in the online version, at
587 XXX.
588 **Appendix B. Model implementation**
589 Model implementation can be found in the online version, at XXX.
590 **Appendix C. Model diagnostics**
591 Model diagnostics can be found in the online version, at XXX.
592 **Appendix D. Simulation study**
593 Simulation study methods and results can be found in the online version, at XXX.
594 **Appendix E. Comparison with two tree-based methods**
595 Two tree-based methods and their results can be found in the online version, at XXX.
596 **Appendix F. Codes for Bayesian variable selection, and tree based methods**
597 Code for Bayesian variable selection and tree-based methods can be found online, at
598 <https://doi.org/10.6084/m9.figshare.5177206>.

599

References

- 600 Albanese, B., Angermeier, P.L., Dorai-Raj, S., 2004. Ecological correlates of fish
601 movement in a network of Virginia streams. *Can. J. Fish. Aquat. Sci.* 61, 857–869.
602 doi:10.1139/f04-096
- 603 Barbieri, M.M., Berger, J.O., 2004. Optimal predictive model selection. *Ann. Stat.* 32,
604 870–897. doi:10.1214/009053604000000238
- 605 Barbiero, R.P., Schmude, K., Lesht, B.M., Riseng, C.M., Warren, G.J., Tuchman, M.L.,
606 2011. Trends in Diporeia populations across the Laurentian Great Lakes, 1997–
607 2009. *J. Great Lakes Res.* 37, 9–17. doi:10.1016/j.jglr.2010.11.009
- 608 Breiman, L., 2001. Random forests. *Mach. Learn.* 45, 5–32.
609 doi:10.1023/A:1010933404324
- 610 Chezik, K.A., Lester, N.P., Venturelli, P.A., Tierney, K., 2014. Fish growth and degree-
611 days II: selecting a base temperature for an among-population study. *Can. J. Fish.*
612 *Aquat. Sci.* 71, 1303–1311. doi:10.1139/cjfas-2013-0615
- 613 Claeskens, G., Hjort, N.L., 2008. Model selection and model averaging. Cambridge
614 University Press.
- 615 Drouineau, H., Bau, F., Alric, A., Deligne, N., Gomes, P., Sagnes, P., 2017. Silver eel
616 downstream migration in fragmented rivers: use of a Bayesian model to track
617 movements triggering and duration. *Aquat. Living Resour.* 30.
618 doi:10.1051/alr/2017003
- 619 Ebener, M.P., Brenden, T.O., Jones, M.L., 2010a. Estimates of fishing and natural
620 mortality rates for four Lake Whitefish stocks in Northern Lakes Huron and
621 Michigan. *J. Great Lakes Res.* 36, 110–120. doi:10.1016/j.jglr.2009.06.003
- 622 Ebener, M.P., Brenden, T.O., Wright, G.M., Jones, M.L., Faisal, M., 2010b. Spatial and
623 temporal distributions of lake whitefish spawning stocks in Northern lakes Michigan
624 and Huron, 2003–2008. *J. Great Lakes Res.* 36, 38–51.
625 doi:10.1016/j.jglr.2010.02.002
- 626 Ethier, D.M., Koper, N., Nudds, T.D., 2017. Spatiotemporal variation in mechanisms
627 driving regional-scale population dynamics of a Threatened grassland bird. *Ecol.*
628 *Evol.* 7, 4152–4162. doi:10.1002/ece3.3004
- 629 Forseth, T., Nesje, T.F., Jonsson, B., Harsaker, K., 1999. Juvenile migration in brown
630 trout: a consequence of energetic state. *J. Anim. Ecol.* 68, 783–793.
631 doi:10.1046/j.1365-2656.1999.00329.x
- 632 Fu, C., Fanning, L.P., 2004. Spatial Considerations in the Management of Atlantic Cod
633 off Nova Scotia, Canada. *North Am. J. Fish. Manag.* 24, 775–784.
634 doi:10.1577/M03-134.1
- 635 Gatz, A.J., Adams, S.M., 1994. Patterns of movement of centrarchids in two warmwater
636 streams in eastern Tennessee. *Ecol. Freshw. Fish* 3, 35–48. doi:10.1111/j.1600-
637 0633.1994.tb00105.x
- 638 Gelman, A., Meng, X.-L., Stern, H., 1996. Posterior predictive assessment of model
639 fitness via realized discrepancies. *Stat. Sin.* 6, 733–807. doi:10.1.1.142.9951
- 640 Gilliam, J.F., Fraser, D.F., 2001. MOVEMENT IN CORRIDORS: ENHANCEMENT
641 BY PREDATION THREAT, DISTURBANCE, AND HABITAT STRUCTURE.
642 *Ecology* 82, 258–273. doi:10.1890/0012-9658(2001)082[0258:MICEBP]2.0.CO;2
- 643 Green, P.J., 1995. Reversible jump Markov chain Monte Carlo computation and Bayesian
644 model determination. *Biometrika* 82, 711–732. doi:10.1093/biomet/82.4.711

645 Gunning, G.E., Shoop, C.R., 1963. Occupancy of home range by longear sunfish,
646 *Lepomis m. megalotis* (Rafinesque), and bluegill, *Lepomis m. macrochirus*
647 Rafinesque. *Anim. Behav.* 11, 325–330. doi:10.1016/S0003-3472(63)80119-0
648 Hutchings, J.A., 1996. Spatial and temporal variation in the density of northern cod and a
649 review of hypotheses for the stock's collapse. *Can. J. Fish. Aquat. Sci.* 53, 943–962.
650 doi:10.1139/f96-097
651 Li, Y., Bence, J.R., Brenden, T.O., 2015. An evaluation of alternative assessment
652 approaches for intermixing fish populations: a case study with Great Lakes lake
653 whitefish. *ICES J. Mar. Sci.* 72, 70–81. doi:10.1093/icesjms/fsu057
654 Lidicker, W.Z., Stenseth, N.C., 1992. To disperse or not to disperse: who does it and
655 why?, in: *Animal Dispersal*. Springer Netherlands, Dordrecht, pp. 21–36.
656 doi:10.1007/978-94-011-2338-9_2
657 McNickle, G.G., Rennie, M.D., Sprules, W.G., 2006. Changes in Benthic Invertebrate
658 Communities of South Bay, Lake Huron Following Invasion by Zebra Mussels
659 (*Dreissena polymorpha*), and Potential Effects on Lake Whitefish (*Coregonus*
660 *clupeaformis*) Diet and Growth. *J. Great Lakes Res.* 32, 180–193. doi:10.3394/0380-
661 1330(2006)32[180:CIBICO]2.0.CO;2
662 Minns, C.K., 1995. Allometry of home range size in lake and river fishes. *Can. J. Fish.*
663 *Aquat. Sci.* 52, 1499–1508. doi:10.1139/f95-144
664 Mitchell, T.J., Beauchamp, J.J., 1988. Bayesian Variable Selection in Linear Regression.
665 *J. Am. Stat. Assoc.* 83, 1023–1032.
666 Mohr, L.C., Nalepa, T.F., 2005. Proceedings of a workshop on the dynamics of lake
667 whitefish (*Coregonus clupeaformis*) and the amphipod *Diporeia* spp. in the Great
668 Lakes, Great Lakes Fish. Comm. Tech. Rep.
669 Polacheck, T., Eveson, J.P., Laslett, G.M., Pollock, K.H., Hearn, W.S., 2006. Integrating
670 catch-at-age and multiyear tagging data: a combined Brownie and Petersen
671 estimation approach in a fishery context. *Can. J. Fish. Aquat. Sci.* 63, 534–548.
672 doi:10.1139/f05-232
673 Price, H., Pothoven, S.A., McCormick, M.J., Jensen, P.C., Fahnenstiel, G.L., 2003.
674 Temperature Influence on Commercial Lake Whitefish Harvest in Eastern Lake
675 Michigan. *J. Great Lakes Res.* 29, 296–300. doi:10.1016/S0380-1330(03)70434-1
676 Radinger, J., Wolter, C., 2014. Patterns and predictors of fish dispersal in rivers. *Fish*
677 *Fish.* 15, 456–473. doi:10.1111/faf.12028
678 Rättsch, G., Onoda, T., Müller, K.-R., 2001. Soft Margins for AdaBoost. *Mach. Learn.* 42,
679 287–320. doi:10.1023/A:1007618119488
680 Rennie, M.D., Ebener, M.P., Wagner, T., 2012. Can migration mitigate the effects of
681 ecosystem change? Patterns of dispersal, energy acquisition and allocation in Great
682 Lakes lake whitefish (*Coregonus clupeaformis*). *Adv. Limnol.* 63, 455–476.
683 doi:10.1127/advlim/63/2012/455
684 Rennie, M.D., Sprules, W.G., Johnson, T.B., 2009. Factors affecting the growth and
685 condition of lake whitefish (*Coregonus clupeaformis*). *Can. J. Fish. Aquat. Sci.* 66,
686 2096–2108. doi:10.1139/F09-139
687 Roff, D.A., 1991. Life History Consequences of Bioenergetic and Biomechanical
688 Constraints on Migration. *Am. Zool.* 31, 205–216. doi:10.1093/icb/31.1.205
689 Rothschild, B.J., 2007. Coherence of Atlantic Cod Stock Dynamics in the Northwest
690 Atlantic Ocean. *Trans. Am. Fish. Soc.* 136, 858–874. doi:10.1577/T06-213.1

691 Runge, C.A., Martin, T.G., Possingham, H.P., Willis, S.G., Fuller, R.A., 2014.
692 Conserving mobile species. *Front. Ecol. Environ.* 12, 395–402. doi:10.1890/130237
693 Schneider, K.N., Newman, R.M., Card, V., Weisberg, S., Pereira, D.L., 2010. Timing of
694 Walleye Spawning as an Indicator of Climate Change. *Trans. Am. Fish. Soc.* 139,
695 1198–1210. doi:10.1577/T09-129.1
696 Schwarz, G., 1978. Estimating the Dimension of a Model. *Ann. Stat.* 6, 461–464.
697 doi:10.1214/aos/1176344136
698 Vandergoot, C.S., Brenden, T.O., 2014. Spatially Varying Population Demographics and
699 Fishery Characteristics of Lake Erie Walleyes Inferred from a Long-Term Tag
700 Recovery Study. *Trans. Am. Fish. Soc.* 143, 188–204.
701 doi:10.1080/00028487.2013.837095
702 Vincenty, T., 1975. DIRECT AND INVERSE SOLUTIONS OF GEODESICS ON THE
703 ELLIPSOID WITH APPLICATION OF NESTED EQUATIONS. *Surv. Rev.* 23,
704 88–93. doi:10.1179/sre.1975.23.176.88
705 Woznicki, S.A., Nejadhashemi, A.P., Abouali, M., Herman, M.R., Esfahanian, E.,
706 Hamaamin, Y.A., Zhang, Z., 2016. Ecohydrological modeling for large-scale
707 environmental impact assessment. *Sci. Total Environ.* 543, 274–286.
708 doi:10.1016/j.scitotenv.2015.11.044
709

710 Table 1. Summary of candidate variables/terms and their interpretation and relationship
711 to hypotheses. A. For continuous variables, “Hypothesis” (first column) states our a priori
712 hypothesis associated with the variable, and the second column indicates sign of
713 associated coefficient that would support that hypothesis. B. Similarly for interaction
714 terms, but here a single hypothesis (our GDD Hypothesis 2) is associated with all
715 interaction terms, and the second column describes the interpretation of coefficients and
716 the pattern in their sign that would support the hypothesis. C. For categorical (dummy)
717 variables we did not have explicit a priori hypothesis for the sign of coefficients but did
718 hypothesize that these factors could influence net distance. For these variables one level
719 of a factor is the baseline with coefficient fixed at zero, and this level (category) is given
720 in the first column and interpretation of the sign of other coefficients in the second
721 column. For A through C, “X” in the “GDD H1” column indicates that the variable was a
722 candidate variable/term in our variable selection process for the GDD Hypothesis 1 Case,
723 and the GDD H2 column likewise indicates if the variable/term was a candidate variable
724 for the GDD Hypothesis 2 Case. The last row summarizes the total number of candidate
725 variables for each GDD hypothesis.

726

727 Table 2. Posterior mean and 95% credible intervals for parameters of the highest
728 posterior probability model.

729

730 Fig. 1. Map of the study area (Lake Huron) and seven tag release (spawning) sites. Of
731 total 1368 recoveries, 659 were from Detour, 300 from Cheboygan, 243 from Burnt
732 Island, 42 from Saginaw Bay, 43 from Sarnia, 56 from Alpena, and 25 from Fishing
733 Islands.

734

735 Fig. 2. Posterior distributions for the number of selected variables (i.e., $q - 1$). The x-axis
736 starts at 5 because all models selected at least five variables.

737

738 Fig. 3. Variable-wise summary results (posterior mean with 95% credible intervals) of
739 the effect of variables (the β_j), with variables named on y-axis for the case with GDD
740 hypothesis 1. Bars are highlighted by red color when the 95% credible interval does not
741 cover 0, which is defined as a consistent effect. The number above each bar is the
742 marginal inclusion probability.

743

744 Fig. 4. Variable-wise summary results (posterior mean with 95% credible interval of the
745 effect β_j for the j th variable, as indicated in y-axis) for the case with GDD hypothesis 2.
746 Bars are highlighted by red color when the 95% credible interval does not cover 0, which
747 was defined as a consistent effect. The number above each bar is the marginal inclusion
748 probability.

749

750 Fig. 5. Model-wise summary for top 12 models ranked according to their posterior
751 probability mass, for the case of GDD hypothesis 2. Variables that were included in the
752 top 12 models are given on the y-axis. Horizontal bar represents posterior 95% credible
753 intervals and symbols on each bar the posterior mean for each coefficient included in a
754 model, with the associated model given to the left of the bar. Thus when more bars are
755 given for a variable it was included in more models.

Table 1

A. Continuous Variables				
Variable Name	Hypothesis	If support, sign of covariate	GDD H1	GDD H2
Length	Greater total length, fish range further from tagging site.	>0	X	X
years_lag	Longer lag between tagging and recovery year, recoveries tend to be further from tagging site.	>0	X	X
<i>Diporeia</i>	Higher relative <i>Diporeia</i> spp. density near the tagging site, fish stay closer to their tagging site.	<0	X	X
GDD_Diff	Greater GDD at the tagging location than the lake average, fish stay closer to their tagging site.	<0	X	
B. Interaction Terms				
Names	Hypothesis	Sign of coefficient	GDD H1	GDD H2
Sep× GDD _{lake} Oct× GDD _{lake} Nov× GDD _{lake} Dec× GDD _{lake}	In years when lake average GDD is higher there is a shift in spawning timing. This is reflected in shorter net distances in one or more adjacent spawning months, and longer net distances in later months.	If <0, fish are closer to tagging site with higher GDD _{lake} during that month, and if >0 further away. Support for hypothesis would be >0 coefficient for one or more adjacent months of Sep – Nov, and <0 coefficient for later months.		X
C. Categorical (Dummy) Variables				
Variable Names	Baseline category (effect was 0)	Interpretation of coefficient	GDD H1	GDD H2
tag_site: Cheboygan tag_site: Burnt_Island tag_site: Alpena tag_site: Sarnia	Fish tagged and released from Detour (Figure 1)	If >0, larger net distance than baseline; if <0, shorter net distance than baseline	X	X

tag_site:				
Saginaw_Bay				
tag_site:				
Fish_Islands				
sex: Female	Male tagged fish	Same as above	X	X
rec_Y: 2003	Tagged fish recovered	Same as above	X	X
rec_Y: 2004	from 2006			
rec_Y: 2005				
rec_Y: 2007				
rec_Y: 2008				
rec_Y: 2009				
rec_Y: 2010				
rec_Y: 2011				
rec_M:7	Tagged fish recovered in	Same as above	X	X
rec_M:8	June of each year			
rec_M:9				
rec_M:10				
rec_M:11				
rec_M:12				
rec_M:1				
rec_M:2				
rec_M:3				
rec_M:4				
rec_M:5				
tag_Y: 2003	Fish tagged and released	Same as above	X	X
tag_Y: 2005	from 2004			
tag_Y: 2006				
Total number of candidate variables for each case (without intercept)			33	36

Table 2

Variable	Mean	Lower	Upper
Rec_M:11	-0.49	-0.63	-0.35
Oct \times GDD _{lake}	-0.45	-0.58	-0.32
<i>Diporeia</i>	-0.17	-0.22	-0.12
length	0.09	0.05	0.14
tag_site: Cheboygan	0.69	0.57	0.80
tag_site: Alpena	1.04	0.78	1.30

Figure 1

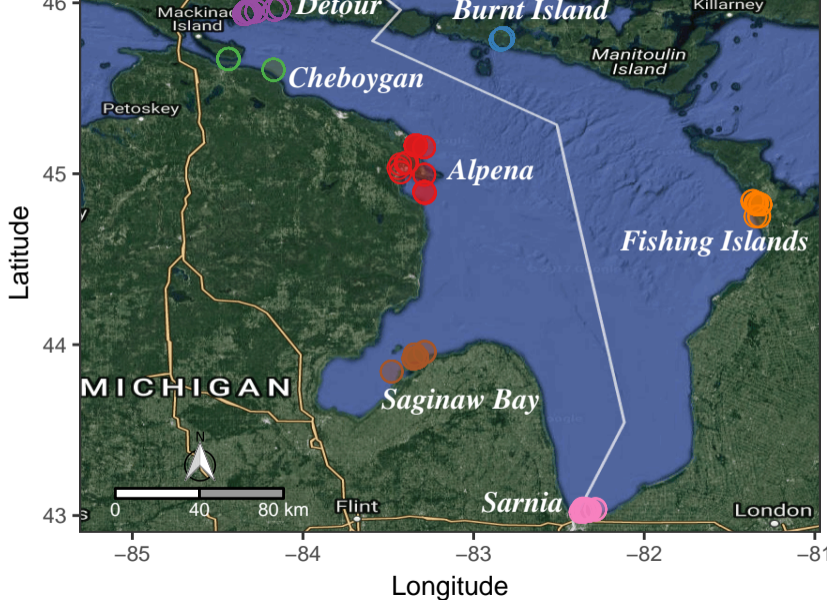


Figure 2

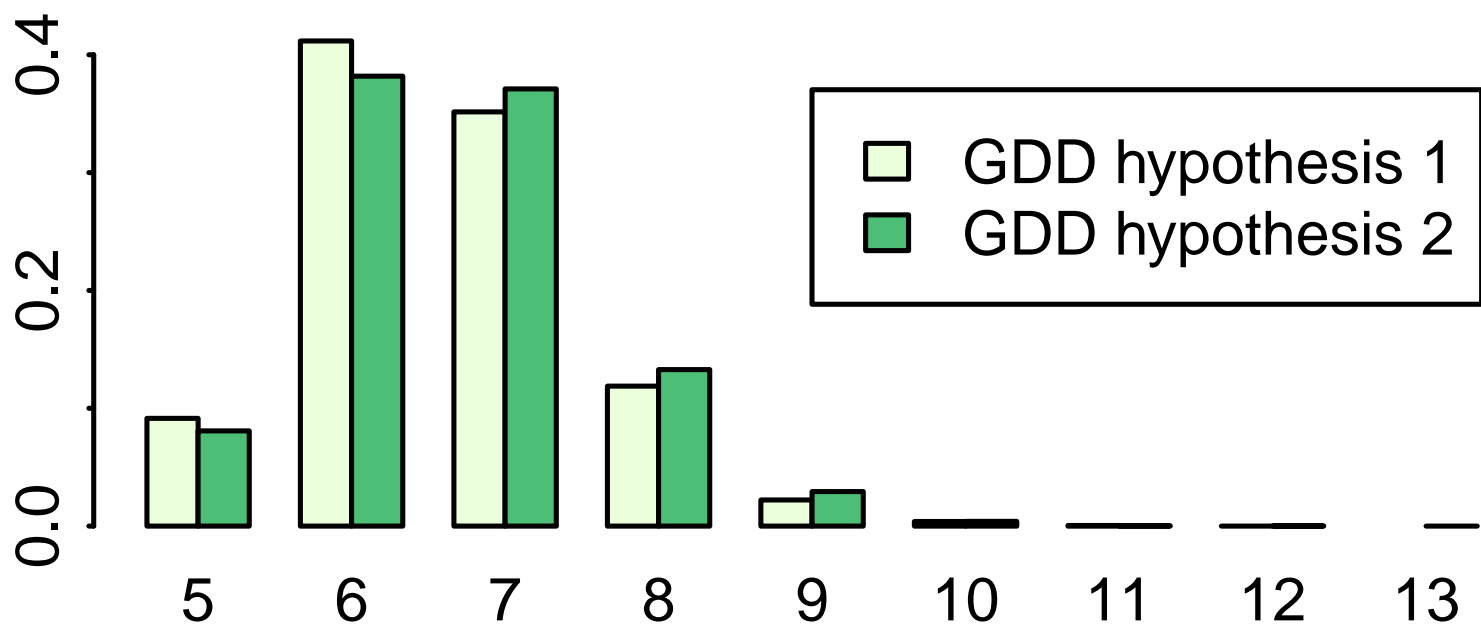


Figure 3

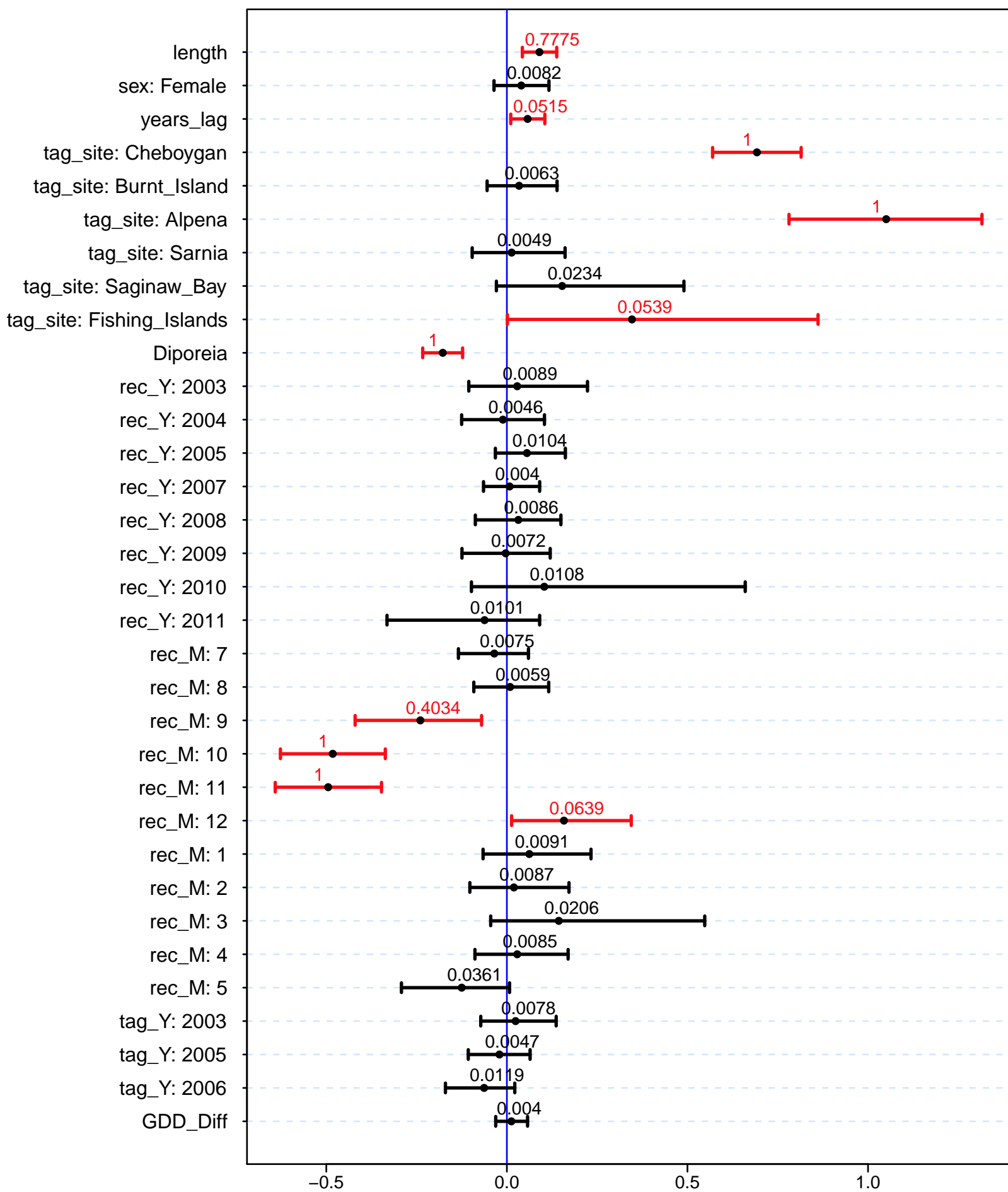


Figure 4

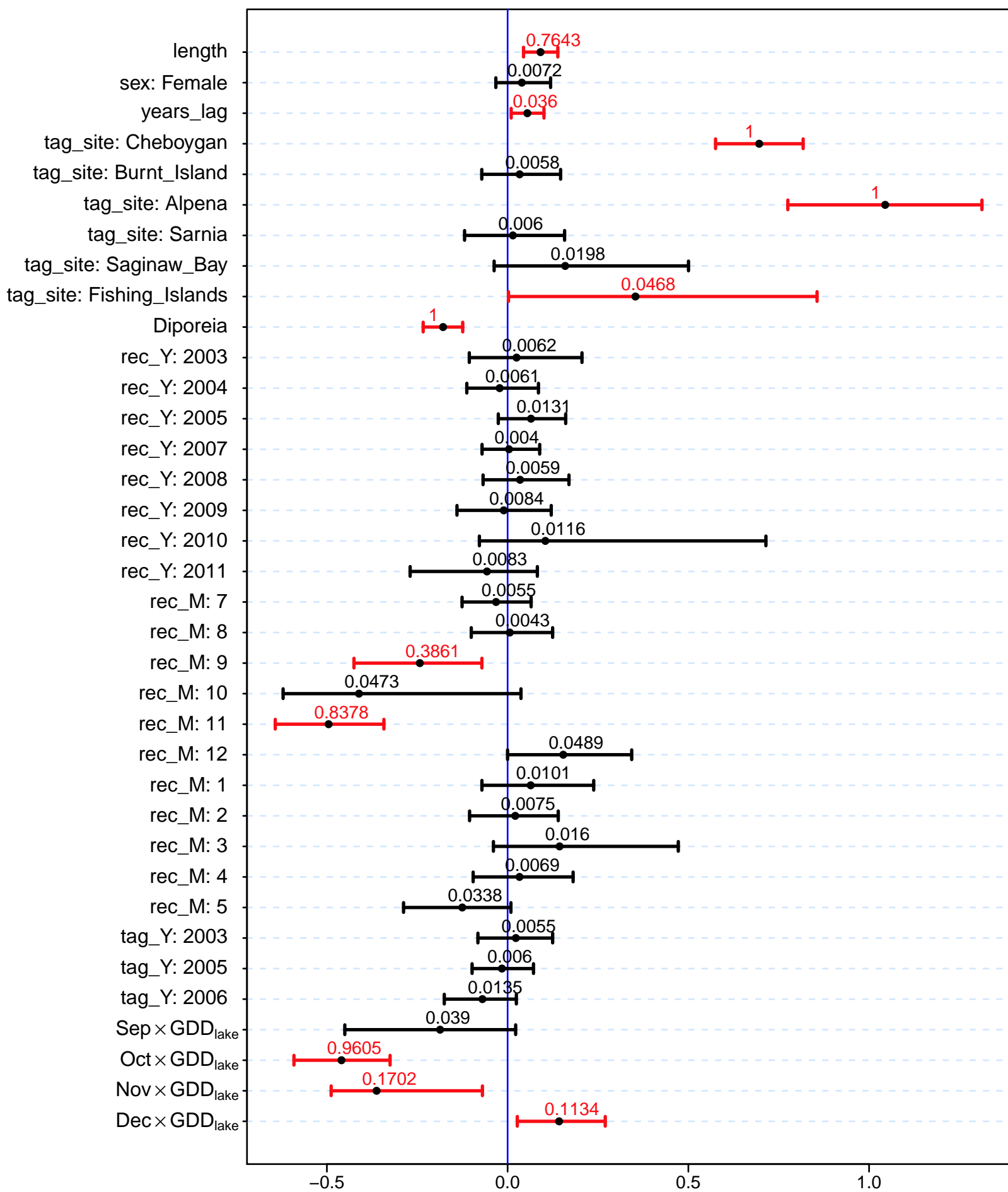
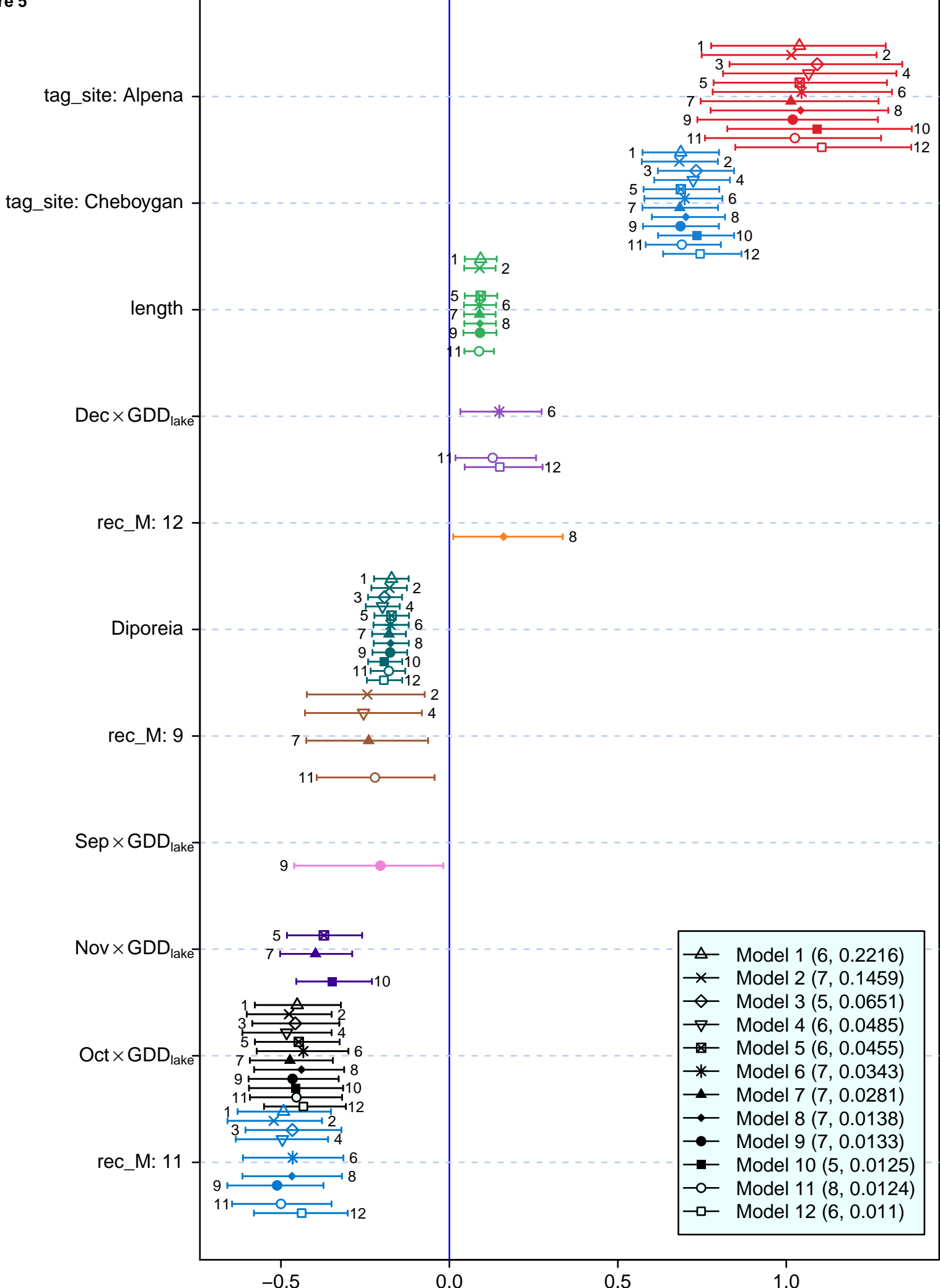


Figure 5



Supplementary material for on-line publication only

[Click here to download Supplementary material for on-line publication only: Second Revision_Supplementary Materials.docx](#)