

1
2 **A skill assessment framework for the Fisheries and Marine Ecosystem Model**
3 **Intercomparison Project**
4 **Nina Rynne^{1,2*}, Camilla Novaglio^{2,3}, Julia Blanchard^{2,3}, Daniele Bianchi⁴, Villy**
5 **Christensen^{5,6}, Marta Coll^{6,7}, Jerome Guiet⁴, Jeroen Steenbeek⁶, Andrea Bryndum-**
6 **Buchholz⁸, Tyler D. Eddy⁸, Cheryl Harrison⁹, Olivier Maury¹⁰, Kelly Ortega-Cisneros¹¹,**
7 **Colleen M. Petrik¹², Derek P. Tittensor¹³ & Ryan F. Heneghan^{1,14, 15}**

8 ¹ School of Mathematical Sciences, Queensland University of Technology, 4 George St,
9 Brisbane, Queensland, Australia.

10 ² Institute for Marine and Antarctic Studies, University of Tasmania, Castray Esplanade,
11 Hobart, Tasmania, Australia.

12 ³ Centre for Marine Socioecology, University of Tasmania, Castray Esplanade, Hobart,
13 Tasmania, Australia.

14 ⁴ Department of Atmospheric and Oceanic Sciences, University of California Los Angeles,
15 Los Angeles, CA, USA.

16 ⁵ Institute for the Oceans and Fisheries, University of British Columbia, Vancouver, Canada.

17 ⁶ Ecopath International Initiative, Barcelona, Spain.

18 ⁷ Institute of Marine Sciences, Barcelona, Spain.

19 ⁸ Centre for Fisheries Ecosystem Research, Fisheries & Marine Institute, Memorial
20 University, St John's, NL, Canada.

21 ⁹. Department of Ocean and Coastal Science and Center for Computation and Technology,
22 Louisiana State University, Baton Rouge, Louisiana, USA.

23 ¹⁰. IRD, Univ. Montpellier, Ifremer, CNRS, INRAE, MARBEC, Sète France.

24 ¹¹. Department of Biological Sciences, University of Cape Town, Cape Town, South Africa.

25 ¹². Scripps Institution of Oceanography, University of California San Diego, La Jolla, CA, USA.

26 ¹³. Department of Biology, Dalhousie University, Halifax, Canada.

27 ¹⁴. School of Science, Technology and Engineering, University of the Sunshine Coast, Petrie,
28 Queensland, Australia.

29 ¹⁵. Australian Rivers Institute, School of Environment and Science, Griffith University,
30 Nathan, QLD, Australia

31 * Corresponding author: Nina Rynne (nina.rynnne@qut.edu.au)

32

33 **Key Points:**

- 34 • We developed a standardised skill assessment framework for an ensemble of global
35 marine ecosystem models
- 36 • Selected models show agreement with the trajectory of fisheries catch, but exhibit
37 biases compared to observed absolute catch values
- 38 • Our framework provides a solid basis to guide global marine ensemble model
39 improvement and increase credibility of ensemble projections

40

41 Abstract

42 Understanding climate change impacts on global marine ecosystems and fisheries requires
43 complex marine ecosystem models, forced by global climate projections, that can robustly
44 detect and project changes. The Fisheries and Marine Ecosystems Model Intercomparison
45 Project (FishMIP) uses an ensemble modelling approach to fill this crucial gap. Yet FishMIP
46 does not have a standardised skill assessment framework to quantify the ability of member
47 models to reproduce past observations and to guide model improvement. In this study, we
48 apply a comprehensive model skill assessment framework to a subset of global FishMIP
49 models that produce historical fisheries catches. We consider a suite of metrics and assess
50 their utility in illustrating the models' ability to reproduce observed fisheries catches. Our
51 findings reveal improvement in model performance at both global and regional (Large
52 Marine Ecosystem) scales from the Coupled Model Intercomparison Project Phase 5 and 6
53 simulation rounds. Our analysis underscores the importance of employing easily
54 interpretable, relative skill metrics to estimate the capability of models to capture temporal
55 variations, alongside absolute error measures to characterise shifts in the magnitude of
56 these variations between models and across simulation rounds. The skill assessment
57 framework developed and tested here provides a first objective assessment and a baseline
58 of the FishMIP ensemble's skill in reproducing historical catch at the global and regional
59 scale. This assessment can be further improved and systematically applied to test the
60 reliability of FishMIP models across the whole model ensemble from future simulation
61 rounds and include more variables like fish biomass or production.

62 **1 Introduction**

63 Across the world's oceans, marine ecosystems are impacted by humans through fishing,
64 pollution, land use change, and via the accelerating impacts of climate change and
65 ecosystem degradation (Halpern et al., 2008; Hatton et al., 2021). Demand on marine
66 ecosystems for food production is already outpacing human population growth (FAO,
67 2022), while climate change impacts are expected to perturb marine communities, from
68 individuals to ecosystems (Fulton et al., 2019), driving changes in the availability,
69 resilience, biomass and location of fish stocks (Blanchard et al., 2012; Booth et al., 2017;
70 Cheung et al., 2010; Hollowed et al., 2013; Lotze et al., 2019; Tittensor et al., 2021).

71 With the growing scope of human impacts on life below water, a range of global Marine
72 Ecosystem Models (MEMs) has been developed by the international research community to
73 help understand and project future change; from simple models based on macroecological
74 scalings to end-to-end models that explicitly represent physical, ecological, and human
75 dynamics; spanning regional systems up to the global ocean. The Fisheries and Marine
76 Ecosystem Model Intercomparison Project (FishMIP; Lotze et al., 2019; Tittensor et al.,
77 2018, 2021; www.fishmip.org) was established in 2013 as part of the broader Inter-
78 Sectoral Impact Model Intercomparison Project (ISIMIP; www.isimip.org) to capitalise on
79 the benefits of bringing such models together into an ensemble. As of today, the FishMIP
80 ensemble comprises no less than nine global and over thirty regional MEMs (Tittensor et
81 al., 2021; Ortega-Cisneros et al., this issue). Individual MEMs are forced by standardised
82 inputs to investigate the influence of environmental conditions and global fishing on ocean
83 biomass and catches while accounting for structural uncertainty across the models
84 (Tittensor et al., 2018, 2021).

85 Amongst today's most relevant applications of MEMs is quantifying and projecting
86 anthropogenic impacts on marine ecosystems with an overarching goal of informing
87 climate change mitigation and adaptation policy, food security issues, and biodiversity
88 policy (Novaglio et al., 2024). Yet to credibly project anthropogenic impacts on marine
89 ecosystems, the reliability of these MEM projections must be assessed in terms of their skill
90 in reproducing past observations. Skill assessment, broadly, involves comparing model
91 outputs for each member of an ensemble with independent sources of observational data
92 using statistical techniques, and comparing metrics of skill across models (Baumberger et
93 al., 2017; Kubicek et al., 2015; Power, 1993; Stow et al., 2009). Robust skill assessment for
94 ecological models is challenging and still lacks widespread usage. A recent review found
95 that as few as 24% of published ecological modelling studies conducted some form of
96 objective (i.e. metric based) skill assessment (Kubicek et al., 2015), while basic visual
97 comparison is the most commonly used subjective skill assessment method, and is
98 arguably the de facto community standard (Stow et al., 2009). More work is needed to
99 standardise and accelerate the use of model skill assessments to enhance the credibility
100 and reliability of ecosystem model projections for MEMs.

101 Rigorous model skill assessment needs to address the relevance of models to the scientific
102 or societal question they are addressing (Jakeman et al., 2006; Planque et al., 2022). This
103 includes identifying sets of relevant metrics that help quantify the realism of simulated
104 variables or patterns of importance (Allen & Somerfield, 2009; Bennett et al., 2013; Power,
105 1993; Stow et al., 2009); and addressing the relevance of the models given technical
106 limitations or the needs of end-users (Hamilton et al., 2019; Kubicek et al., 2015; Steenbeek
107 et al., 2021, this issue).

108 There are four major challenges that have inhibited the widespread usage of skill
109 assessment for ecological models including MEMs, and especially for cross-model
110 comparison: (1) the absence of a standardised framework, leading to arbitrary and
111 inconsistent choices of important metrics (Geary et al., 2020; Hipsey et al., 2020; Kubicek et
112 al., 2015; Mayer & Butler, 1993; Rykiel, 1996); (2) the need for multiple relevant metrics to
113 assess different aspects of model performance, as relying on a single metric can obscure
114 divergent behaviour or favour models that are highly correlated with a particular set of
115 observations by chance (Bennett et al., 2013; Eyring et al., 2019; Legates & McCabe Jr.,
116 1999; Mayer & Butler, 1993; Power, 1993); (3) credible replication of observations by a
117 model in some regions or for a given time-period does not guarantee performance beyond
118 the calibrated range (Eyring et al., 2019; Hipsey et al., 2020; Hollowed et al., 2013;
119 Refsgaard et al., 2014; Steenbeek et al., 2021; Wagener et al., 2022); (4) the hypothetical
120 nature of future projections makes their comparison to observations unfeasible
121 (Baumberger et al., 2017; Hamilton et al., 2019; Hollowed et al., 2013; Refsgaard et al.,
122 2014).

123 In the context of fisheries and ecosystem models, these challenges are compounded by
124 limitations in the observational data available to validate MEMs, and the quality of available
125 data. While global and regional catch reconstructions exist (e.g. Pauly & Zeller, 2016;
126 Watson & Tidd, 2018), global observations of fish biomass are lacking. Stock assessments,
127 such as the RAM Legacy Stock Assessment Database (Ricard et al., 2012;
128 www.ramlegacy.org) or recent regional standardised synthesis of biomass observations
129 from trawl surveys (Maureaud et al., 2023) are filling this gap, though they remain limited

130 in spatial coverage, and represent snapshots in time that do not capture variability at
131 seasonal, interannual, or longer timescales.

132 Although individual skill assessment of FishMIP models has been performed on individual
133 models (Barrier et al., 2023; Blanchard et al., 2012; Carozza et al., 2016, 2017; Cheung et al.,
134 2011; Christensen et al., 2015; Christensen & Walters, 2004; Heneghan et al., 2021;
135 Jennings & Collingridge, 2015; Maury et al., 2007; Maury, 2010; Novaglio et al., 2022;
136 Ortega-Cisneros et al., 2017; Petrik et al., 2019; Sturludottir et al., 2018), these assessments
137 vary from model to model and no standardised skill assessment across the FishMIP model
138 ensemble has taken place yet.

139 This paper sets the foundations of a skill assessment framework for FishMIP based on the
140 “Concept-State-Process-System” (CSPS) framework created by Hipsey et al. (2020), to build
141 confidence that future predictions are robust and credible. We choose the CSPS framework
142 for its thorough, multi-step process and extensive integration of skill assessment examples
143 in aquatic ecosystem literature.

144 In choosing and adapting the CSPS framework we attempt to address the first three key
145 challenges (absence of a standardised framework, the need for multiple metrics and model
146 credibility beyond calibrated range) that have hindered the widespread use of model skill
147 assessment in marine ecosystem modelling, and FishMIP in particular. In doing so we: (1)
148 conduct the first standardised model skill assessment of fisheries catch predictions across a
149 subset of FishMIP ensemble members and (2) demonstrate how the CSPS framework can
150 be utilised as the foundation for building context-specific ecosystem skill assessment tools.

151 We argue that the subsequent uptake of such a framework will improve the credibility and
152 reliability of global MEMs, strengthening their use to inform decision-making. The wider
153 benefits of this case study include illustrating how the framework can be customised for
154 the skill assessment needs of other model intercomparison projects, and therefore further
155 the development of open-access, reliable skill assessment tools for other climate-impact
156 assessment ensembles.

157 **2 Materials and Methods**

158 **2.1 Concept-State-Process-System (CSPS) Framework Overview**

159 The CSPS framework categorises measurable elements of ecosystem structure and function
160 across four levels (Hipsey et al., 2020). Level 0 (conceptual assessment) focuses on
161 ensuring that model parameterisations, assumptions and representation of the underlying
162 system are reasonable and derived from a credible scientific basis, and that the level of
163 complexity in the model is appropriate for the questions being asked (Hipsey et al., 2020;
164 Kubicek et al., 2015; Rykiel, 1996). Level 1 (state assessment) is concerned with a model's
165 ability to reflect past observations of measured ecosystem properties. This is generally
166 achieved using metrics that assess goodness-of-fit, which can highlight various mismatches
167 between observations and simulations (Stow et al., 2009). Level 2 (process assessment)
168 and Level 3 (system assessment) explore whether the model is right for the right reasons,
169 i.e., whether the model has captured the important underlying rates of change within the
170 ecosystem, as well as spatial and temporal dynamics that emerge from the model (Hipsey
171 et al., 2020).

172 Due to the complex task of developing a standardised skill assessment framework for
173 FishMIP models (as detailed in the introduction) and the current lack of a set of
174 quantitative measures to assess FishMIP models' ability to reproduce past trends and
175 patterns, this paper focuses on the implementation of Level 1 model skill assessment.
176 However, we will consider future developments, including the development of methods for
177 Levels 2 and 3 in the context of FishMIP in the Discussion, noting that such methods are
178 still broadly described in the literature and seldom considered when assessing marine
179 ecosystem models.

180 CSPS addresses three of the four major challenges of MEMs skill assessment: it provides a
181 standardised approach and selection of metrics that can be used across models (addressing
182 the absence of a standardised assessment framework); multiple metrics are used to assess
183 model performance (addressing the need for a suite of metrics to holistically assess model
184 performance); and, finally it provides a framework to assess whether models can replicate
185 ecosystem processes and properties, which is critical for MEMs to provide credible
186 prediction outside their calibrated range (Hipsey et al., 2020; Kubicek et al., 2015;
187 Steenbeek et al., 2021). The identification of emergent processes are context- and model-
188 specific, as the important dynamics to be assessed vary depending on the purpose of the
189 model (Petrik et al., 2022; Novaglio et al., this issue). Identifying relationships between
190 historical simulations and forecasted ocean states through emergent constraints will help
191 address the challenge arising from the hypothetical nature of model projections.

192 **2.2 Metrics for Level 1 FishMIP model assessment**

193 We adapted the CSPA framework to the skill assessment needs of FishMIP in three steps.

194 First, we used the CSPA framework as a benchmark to categorise model skill assessment

195 approaches proposed in other papers (see Table S1). Second, we used best practices and

196 commonly agreed-upon statistical measures to recommend FishMIP-appropriate

197 assessment measures (Table 1). Third, we applied the framework to two structurally

198 contrasting global FishMIP models and assessed their respective ability to reproduce

199 historical fisheries catches at both the global and Large Marine Ecosystem (LME) scale.

200 We used a synthesis of goodness-of-fit metrics from existing skill assessment literature

201 (Table 1, S2), from which we selected a range of statistical measures for this study (Table

202 1). Following the advice of Legates & McCabe Jr. (1999) in using both relative and absolute

203 error measures, for our Level 1 assessment we used a visual representation of correlation

204 and bias, and a Taylor diagram plot of standard deviation, Pearson correlation and centred

205 root mean squared error. Alongside this, we calculated 6 independent metrics of model

206 skill assessment (Table 1; Allen & Somerfield, 2009; Bennett et al., 2013; Mayer & Butler,

207 1993; Stow et al., 2009; Taylor, 2001). These metrics measure: (1) the models' ability to

208 replicate trends over time (Pearson correlation (R)); (2) bias between projections and

209 observations (average error (AE), root mean squared error (RMSE), mean absolute error

210 (MAE)); and (3) a combination of trend and bias (reliability index (RI), and modelling

211 efficiency (MEF)). These skill metrics are detailed in Table 1 and S2. These metrics are

212 calculated for the two FishMIP MEMs, considered here and described in the next section,

213 and are tabulated for comparison.

215 **Table 1. Skill Assessment Metrics.** List of skill assessment metrics, their usage, and
 216 additional notes. Adapted from Stow et al., (2009). See Table S2 for more information about
 217 these metrics.

Name	Type	Ideal Value	Usage	Notes
Correlation (R)	Relative	1	R measures the degree to which simulated and observed catches change together in time. This metric indicates if both variables move in the same direction over time.	R is a relative (or dimensionless) statistic, meaning that correlation is not influenced by the magnitude of the underlying data, therefore values close to 1 can occur even if there is considerable difference in magnitude between the values. Additionally, correlation can be sensitive to outliers if they exist in the data. Relative statistics are comparable across different models or regions.
Average Error (AE)	Absolute	0	AE is the sum of the size of the discrepancies between simulated and observed catch value-pairs. It measures the aggregate bias (or under/overestimation) of simulated catches compared to observations.	A shortcoming of AE is that results close to zero can indicate either a close match or can be a result of positive and negative errors cancelling out. To overcome this, other methods of calculating error can be used instead of, or alongside, AE.
Root Mean Squared Error (RMSE)	Absolute	0	RMSE gives the average distance between predicted and observed catches. It measures the aggregate bias of simulated catches compared to observation Centred RMSE – as reported in a Taylor diagram – is given as a RMSE relative to the standard deviation of observed catches.	RMSE accommodates for the shortcomings of AE as it considers the magnitude, but not the direction, of each discrepancy. As RMSE uses the square of each discrepancy, it is more sensitive to the influence of outliers than either AE or MAE.
Mean Absolute Error (MAE)	Absolute	0	MAE is the sum of the absolute size of the discrepancies between simulated and observed catches. It measures the aggregate bias of the simulated catches compared to observations.	MAE accommodates for the shortcoming of AE, by using the absolute value of the discrepancies. When absolute differences are of a similar magnitude, RMSE and MAE will be approximately equal (Mayer & Butler, 1993)
Reliability Index (RI)	Relative	1	RI is a measure of the average multiplicative factor by which simulated catches differ from observations. Similar to AE, RMSE and MAE, it can be used to measure the bias of the simulations, but as a	RI results are relative making them useful for comparing projections from different models or for different regions

relative statistic, it can be compared across different models or regions.

Modelling Efficiency (MEF)	Relative	1	MEF measures the predictive ability of model simulations, relative to the average of the observations in an easily interpretable single statistic.	MEF \in $(-\infty, 1]$. A negative result indicates that the observational average is a better predictor than the model projections. Results >0 indicate that the model is a better predictor than the average of the observations.
----------------------------	----------	---	--	--

218

219 All analyses were carried out in R-statistics version 4.2.3 (R Core Team, 2023). Statistical
 220 metrics were calculated using the R packages “*Metrics*” (Hamner et al., 2018), “*topmodel*”
 221 (Buytaert, 2022) and “*qualV*” (Jachner et al., 2007), and the Taylor diagram was plotted
 222 using the R package “*openair*” (Carslaw & Ropkins, 2012).

223 **2.3 Marine Ecosystem Models**

224 We obtained model data from two published global MEMs that are members of the FishMIP
 225 ensemble and have provided historical simulation outputs of fisheries catches under two
 226 FishMIP simulation protocols (using Coupled Model Intercomparison Project (CMIP) Phase
 227 5 and 6 Earth system model forcings, respectively); the BiOeconomic mArine Trophic Size-
 228 spectrum model (BOATS; Carozza et al., 2016, 2017) and EcoOcean (Christensen et al.,
 229 2015; Coll et al., 2020). These models are a subset of the 9 FishMIP global MEMs (Tittensor
 230 et al., 2021), but they are the only two that had historical outputs of fisheries catches across
 231 the two simulation rounds at time of publication. Nevertheless, these models capture a
 232 significant portion of the spectrum of model complexity across the full FishMIP ensemble,
 233 from BOATS which resolves individual organisms by body size alone, to EcoOcean which
 234 explicitly incorporates information about thousands of species.

235 BOATS is a size-structured model that uses broad-scale ecological relationships and
236 individual-level metabolic constraints to calculate the production of fish biomass. It is
237 coupled with an economic module that determines fishing effort and harvest based on the
238 profitability of the exploitation of this biomass given globally homogenous economic
239 boundary conditions (Carozza et al., 2016, 2017). BOATS fish biomass production is driven
240 by water temperature (averaged over the top 75m) and depth-integrated net primary
241 production from Earth System Models (ESMs) and implicitly includes all commercially
242 fished animal biomass from 10 g to 100 kg for three fish groups of increasing asymptotic
243 mass (0.3 kg, 8.5 kg and 100 kg, respectively). By default, in each grid cell of the simulated
244 domain BOATS assumes open-access fishing effort dynamics (Carozza et al., 2016;
245 Tittensor et al., 2018), although it can be forced by observational reconstructions of effort
246 and other social or economic drivers (Scherrer & Galbraith, 2020).

247 EcoOcean is a combined trophodynamic and species distribution model with a mass-
248 balanced food web model at its core (Christensen et al., 2015; Coll et al., 2020). It uses
249 fishing effort and gear type as forcings, and a gravity model to spatially spread effort across
250 grid cells within LMEs based on expected profitability and fishing costs. Fish prices, used to
251 estimate expected revenue, are model inputs while fishing costs are assumed to be
252 proportional to the grid cell's distance from the nearest coast. Both fish prices and costs are
253 used to calculate fishing effort. The EcoOcean realisation used in this study considers
254 depth-resolved water temperature and depth-integrated small and large phytoplankton
255 biomass as drivers. EcoOcean resolves 51 functional groups, including fish, sharks and rays,
256 invertebrates and mammals, to represent the whole spectrum of marine organisms and

257 integrates explicit information for 3,400 species of marine organisms (Christensen et al.,
258 2015; Coll et al., 2020; Tittensor et al., 2021).

259 Both MEMs use common simulation protocols, as defined by the Inter-Sectoral Impact
260 Model Intercomparison Project (ISIMIP; www.isimip.org). We used output from protocols
261 ISIMIP2b and ISIMIP3b (Blanchard et al., 2024; Frieler et al., 2017, 2024; Tittensor et al.,
262 2018, 2021), which used climate forcings from the CMIP5 and CMIP6, respectively.

263 Total catch output data were provided in a standardised 1° grid cell format monthly.
264 Historical simulations from both MEMs spanned 1971-2005 for ISIMIP2b and 1950-2014
265 for ISIMIP3b. All outputs from the FishMIP ensemble are available at
266 www.isimip.org/gettingstarted/data-access/. This includes outputs used for the two
267 models presented in this study, except for EcoOcean outputs from ISIMIP3b (available
268 here: 10.5281/zenodo.11081600).

269 **3 Marine Ecosystem Model Forcings and Observational Data**

270 For both ISIMIP2b and 3b, BOATS and EcoOcean were forced with outputs from two Earth
271 System Models (ESMs): Geophysical Fluid Dynamics Laboratories (GFDL) (version ESM2M
272 and ESM4.1 for ISIMIP2b and ISIMIP3b, respectively; Dunne et al., 2012, 2020) and Institut
273 Pierre-Simon Laplace (IPSL) (version CM5A-LR and CM6A-LR for ISIMIP2b and ISIMIP3b,
274 respectively; Boucher et al., 2020; Sepulchre et al., 2020). These ESM simulations forced the
275 FishMIP ensemble models for both ISIMIP simulation rounds (e.g., CMIP5 in Lotze et al.,
276 2019; CMIP5 and CMIP6 in Tittensor et al., 2021).

277 Two global fishing catch datasets were initially considered to capture the variability and
278 biases from different data reconstruction methodologies. The first, from Watson & Tidd

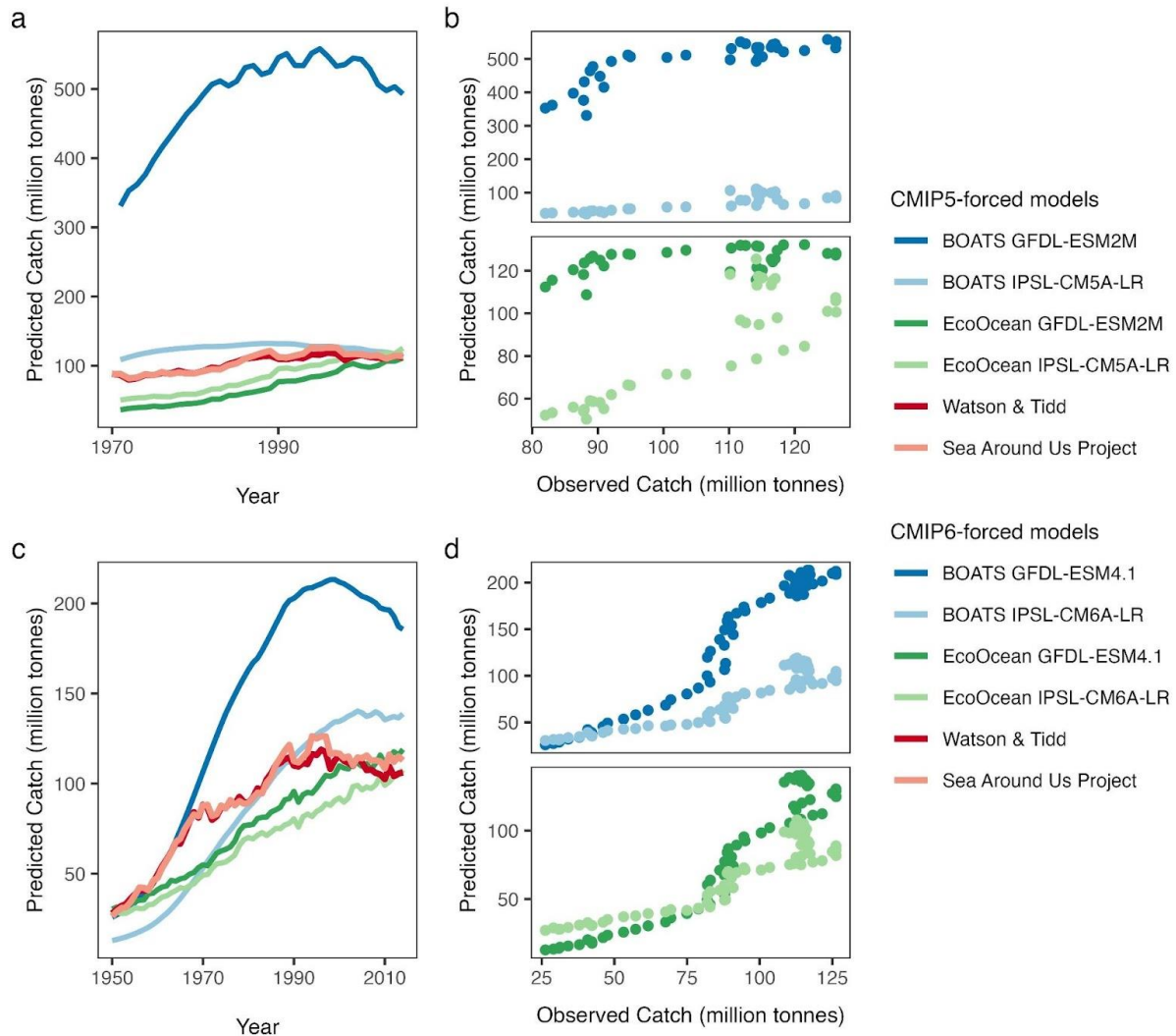
279 (2018), covers the historical period 1869-2017 and combines official reconstructed
280 estimates of fisheries catch data, including major discards, from the Food and Agriculture
281 Organisation (FAO) FishStat database along with other publicly available sources (Watson,
282 2019; Watson & Tidd, 2018; <http://dx.doi.org/10.25959/5c522cadbea37>). The second,
283 from the Sea Around Us Project (SAUP), covers the historical period 1950-2019, uses FAO-
284 reported landings, Regional Fisheries Management Organisations (RFMOs), expert
285 elicitation and other publicly available sources (Pauly & Zeller, 2016;
286 <https://www.seaaroundus.org/data/#/search>). However, at the global scale, the difference
287 between these two datasets is small, relative to the divergence between observations and
288 simulations (Figure 1). Therefore, we used only the Watson & Tidd (2018) reconstruction
289 to calculate model performance metrics.

290 **4 Results**

291 **4.1 Global scale skill assessment**

292 Globally averaged historical time series of simulated catches were strongly correlated with
293 Watson & Tidd observations for both EcoOcean and BOATS (Figure 1b, 1d). CMIP5-forced
294 correlations were lower for both models compared to CMIP6, (Table 2), especially for
295 BOATS-IPSL ($R = 0.47$). For CMIP6-forced simulations, correlation coefficients, R , ranged
296 between 0.92 and 0.98, with slightly higher values for BOATS than for EcoOcean (Table
297 3). Furthermore, bias was substantially lower in CMIP6-forced simulations compared to
298 CMIP5 (Figure 1a vs Figure 1c). Generally, there was greater bias in simulated absolute
299 values of catches compared to observations for BOATS than for EcoOcean across both
300 CMIP5- and 6-forced simulations when using GFDL-forcing (Table 2, Table 3). This result
301 was reversed when using IPSL-forcing, when EcoOcean simulated values showed greater

302 bias. All simulations with EcoOcean generally underestimated catch, while BOATS-GFDL
 303 strongly overestimated catch across CMIP5 and 6, whereas the bias of BOATS-IPSL was
 304 smaller for both CMIP5 and CMIP6 (Figure 1).



305
 306 **Figure 1. Modelled and observed global fishing catch time series.** *Reconstructed*
 307 *observations from Watson & Tidd (2019) and SAUP (2016) and model projected catch for a)*
 308 *CMIP5 from 1971-2005; and c) CMIP6 from 1950-2014. Scatterplot of Watson & Tidd (2019)*
 309 *reconstructed observations vs model predicted catch for b) CMIP5-forced BOATS (top) and*
 310 *EcoOcean (bottom); and d) CMIP6-forced BOATS (top) and EcoOcean (bottom).*

311

312 The calculated errors (AE, MAE and RMSE) confirm large discrepancies in the magnitude of
 313 simulated and observed catches (Table 2, Table 3). In CMIP6-forced models, the negative
 314 result for AE for BOATS-IPSL reflected that simulated catches were lower than
 315 observations before ~1980 and higher after ~1990. This led to an AE result closer to zero
 316 than other MEM simulations because positive and negative measures cancelled each other
 317 out. With the exception of EcoOcean forced by IPSL, all CMIP6-forced models showed
 318 improved results for AE, MAE and RMSE compared to CMIP5-forced models (Table 2, 3).
 319 RMSE results for all CMIP6-forced models were higher than MAE results, potentially
 320 indicating the presence of large outlier values in the simulated catches (Legates & McCabe
 321 Jr., 1999).

322 **Table 2. Global forecast skill metrics for fishing catch with CMIP5 forcing.** Skill metric
 323 performance for six skill metrics using Watson & Tidd (2019) observations: correlation (R), root
 324 mean squared error (RMSE; g/m^2), mean absolute error (MAE; g/m^2), average error (AE; g/m^2),
 325 reliability index (RI) and modelling efficiency (MEF) for BOATS and EcoOcean models under both
 326 ESM forcings. Results in bold are close to ideal results.

MEM-ESM				
	BOATS		EcoOcean	
Skill Metric	GFDL-ESM2M	IPSL-CM5A-LR	GFDL-ESM2M	IPSL-CM5A-LR
R	0.84	0.47	0.83	0.86
RMSE	388,936,857	22,825,452	39,586,313	25,900,241
MAE	385,610,855	19,106,219	36,745,703	22,961,393
AE	385,610,855	19,106,219	-36,745,703	-21,598,979
RI	4.64	1.24	1.71	1.38
MEF	-755.59	-1.61	-6.84	-2.36

327

328

329 **Table 3. Global skill metrics for fishing catch time series with CMIP6 historic forcing.** Skill
 330 metric performance for six skill metrics using Watson & Tidd (2019) observations: correlation (R),
 331 root mean squared error (RMSE; g/m²), mean absolute error (MAE; g/m²), average error (AE; g/m²),
 332 reliability index (RI) and modelling efficiency (MEF) for BOATS and EcoOcean models under both
 333 ESM forcings. Results in bold are close to ideal results.

MEM-ESM				
	BOATS		EcoOcean	
Skill Metric	GFDL-ESM4.1	IPSL-CM6A-LR	GFDL-ESM4.1	IPSL-CM6A-LR
R	0.98	0.95	0.92	0.92
RMSE	65,832,789	20,832,218	17,865,414	26,094,844
MAE	54,295,280	17,874,215	14,633,816	23,209,098
AE	53,625,817	-4,835,888	-13,758,893	-23,181,291
RI	1.57	1.51	1.26	1.41
MEF	-3.99	0.5	0.633	0.216

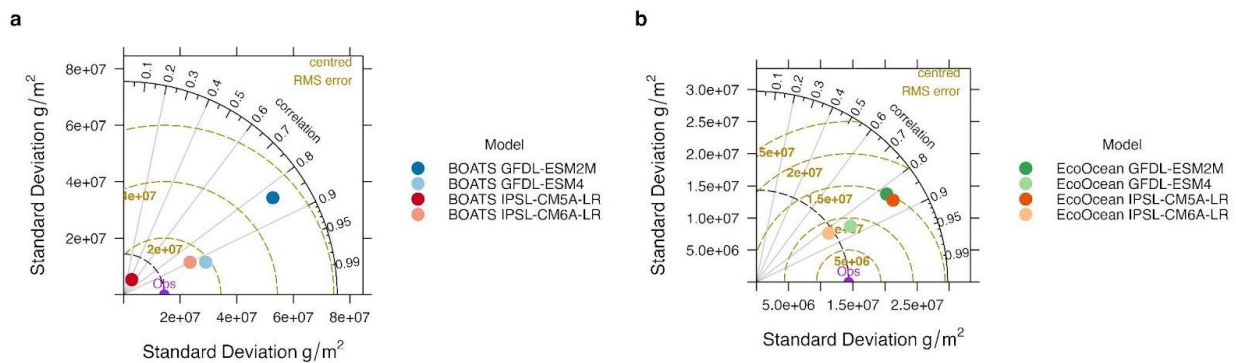
334

335 RI results showed that all CMIP6-forced models differed from fishing catch by between
 336 1.26-1.57-fold on average compared to observations (Table 3). In contrast, for CMIP5
 337 simulated catches, results differed by up to 4.64-fold (BOATS-IPSL) compared to the
 338 observed data (Table 2). This highlighted important improvements in catch estimations
 339 between CMIP5- and 6-forced models. In particular, BOATS_GFDL CMIP6-forced runs
 340 showed the largest improvement, with an RI of 1.57 for BOATS-GFDL, compared to the
 341 CMIP5-forced result of 4.64 (Table 2, 3). Historical catch simulated with CMIP6-IPSL-forced
 342 results worsened slightly compared to CMIP5-IPSL-forced runs (Table 3).

343 For all CMIP5-forced MEM and ESM combinations, modelling efficiency was less than zero,
 344 meaning that the average of the observations is more skillful than the models' estimates
 345 over the historical period (Table 2). In contrast, modelling efficiency (MEF) was greater
 346 than zero for CMIP6-forced BOATS-IPSL, EcoOcean-GFDL and EcoOcean-IPSL (Table 3)
 347 indicating that the simulated catches match observed fishing catches more closely than the
 348 average of the observations for these simulations. However, modelling efficiency remained
 349 negative for CMIP6-forced BOATS-GFDL, albeit greatly improved from CMIP5 (-3.99 versus
 350 -755.59) (Table 3).

351 Taylor diagrams for global catch from BOATS and EcoOcean summarise some of the
 352 previous observations. They show an improvement in correlation and bias between CMIP5-
 353 and CMIP6-forced simulations for both models however, BOATS-IPSL simulations show an
 354 increase in standard deviation from CMIP5- to CMIP6-forced runs (Figure 2).

355



356

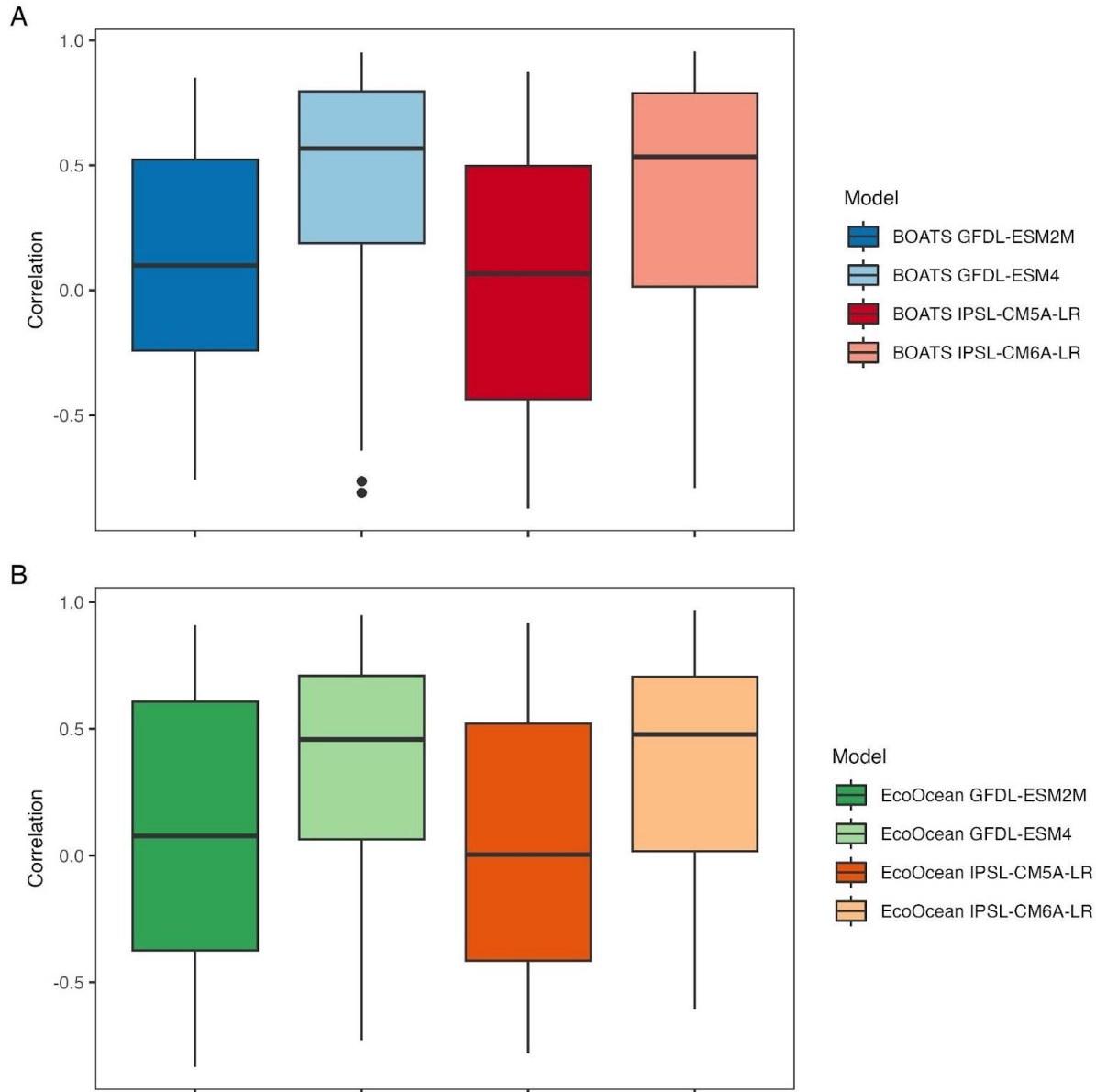
357 **Figure 2. Taylor Diagram for CMIP5 and CMIP6 simulations** a) *BOATS model predicted global*
 358 *catch (from 1971-2005) for CMIP6 using ESM IPSL (light red) and GFDL (light blue), and CMIP5 ESM*
 359 *IPSL (dark red) and GFDL (dark blue); and b) EcoOcean model predicted global catch (from 1971-*
 360 *2005) for CMIP6 using ESM IPSL (light orange) and GFDL (light green), and CMIP5 ESM IPSL (dark*

361 *orange*) and *GFDL* (*dark green*). Plot shows standard deviation, correlation and centred root mean
362 squared error for all 4 MEM-ESM combinations. *Watson & Tidd (2019)* observed global catch
363 (*between 1971-2005*) in purple.

364 **4.2 Large Marine Ecosystem scale assessment**

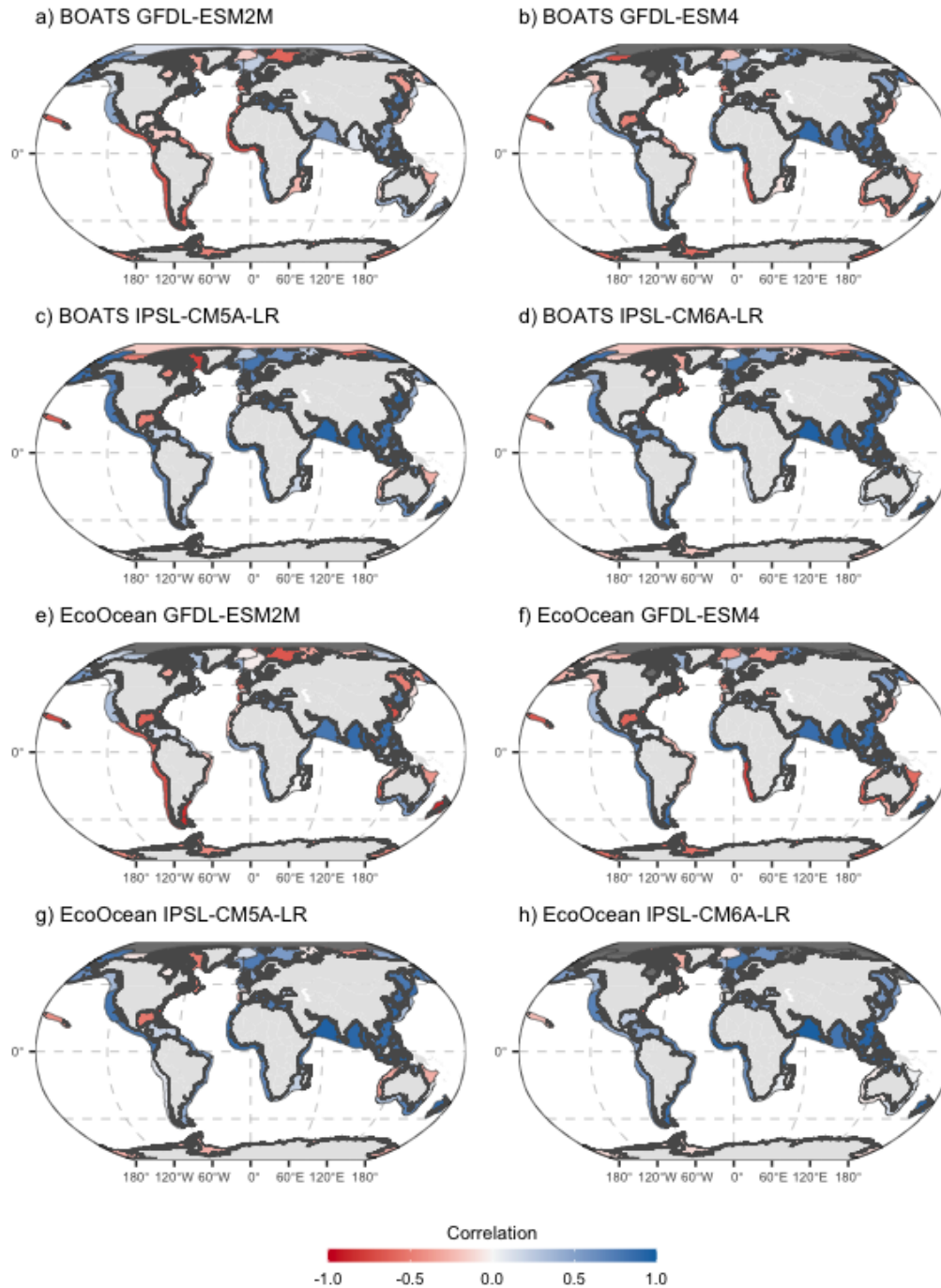
365 Correlations between simulated catches and *Watson & Tidd (2019)* reconstructed
366 observed fishing catch varied across the LMEs, indicating differing levels of model
367 performance at the regional scale. For all CMIP5-forced models, the median correlation
368 (median across LME-levels) was near zero (Figure 3). Results from CMIP6-forced models
369 showed improvement in the median and interquartile range of correlation results
370 compared to CMIP5-forced models (Figure 3), indicating improved correlation at the LME
371 scale overall.

372 Geographically, this improvement in correlation from CMIP5 to CMIP6 was evident for both
373 BOATS and EcoOcean, but the degree of improvement (and where this occurs) differs
374 between ESM-forcings (Figure 4). For BOATS, there was an improvement in correlations in
375 CMIP6 compared to CMIP5 across 50 LMEs with both GFDL and IPSL forcing. Similarly, for
376 EcoOcean there was an improvement in correlations across 47 and 46 LMEs for GFDL and
377 IPSL, respectively (Table S3). BOATS showed marked improvements in highly productive
378 LMEs including in the Humboldt Current, Pacific Central-American Coast, Barents Sea,
379 North Brazil Shelf, Patagonian Shelf and Canary Current (Figure 4; Table S3). In contrast,
380 EcoOcean's largest correlation improvements were more randomly dispersed across
381 European, east African and East South American LMEs (Figure 4; Table S3). Negative
382 correlations between observed and modelled catches persist across all simulations in polar
383 regions (Figure 4).



384

385 **Figure 3. Box plot of Large Marine Ecosystem (LME) level correlation.** a) *BOATS* model
 386 correlation compared to reconstructed observations from Watson & Tidd (2019). CMIP5-forced ESM
 387 IPSL (dark red) and GFDL (dark blue), and CMIP6-forced ESM IPSL (light red) and GFDL (light blue);
 388 b) *EcoOcean* model correlation compared to reconstructed observations from Watson & Tidd
 389 (2019). CMIP5-forced ESM IPSL (dark orange) and GFDL (dark green), and CMIP6-forced ESM IPSL
 390 (light orange) and GFDL (light green).



391

392 **Figure 4. Map of Pearson correlations across the world's LMEs. BOATS and EcoOcean**
 393 *predictions and reconstructed observations from Watson & Tidd (2019), using CMIP5 from 1971-2005*
 394 *(a, c, e, g) and CMIP6 from 1950-2014 (b, d, f, h) ESM forcings.*

395 **5 Discussion**

396 Model skill assessment is essential for improving the credibility and reliability of Marine
397 Ecosystem Model (MEM) simulations and for supporting their use as decision-making tools.
398 Until now Fisheries and Marine Ecosystem Model Intercomparison Project (FishMIP)
399 models have generally been assessed in isolation. Here, we adapted and implemented a
400 standardised skill assessment framework to highlight commonalities and discrepancies
401 between modelled and observed historical fish catch and to investigate the usefulness of a
402 range of skill assessment metrics.

403 Agreement between simulated and observed catches across BOATS and EcoOcean, both in
404 catch time-series variability and absolute catches, is generally higher for CMIP6- than
405 CMIP5-forced models (Figure 1; Table 1, 2). These improvements may be due to changes in
406 the MEMs. For example, EcoOcean has recently undergone substantial restructuring and
407 the upgrades include an expanded food web from 1,400 to over 3,400 explicitly considered
408 individual species, updated functional group representation, and the use of observed
409 historical spatial ranges of species to initialise the model (Coll et al., 2020). This has
410 resulted in an improved understanding and representation of key ecological and fishing
411 dynamics. Between CMIP5 and CMIP6 BOATS' biological formulation and parameters were
412 not changed. However, to improve the match with observed catches, starting effective
413 effort (which then increases through time with improving catchability) was calibrated to
414 align the model's aggregated catch by LMEs with observational reconstructions from Sea
415 Around Us Project (SAUP; Pauly & Zeller, 2016).

416 While these modifications and – in the case of EcoOcean - a reconsideration of model
417 assumptions, in line with Level 0 (conceptual) assessment, are likely substantial drivers of
418 each model's better performance from CMIP5 to CMIP6, some of the improvement in
419 simulated fishing catches would reflect changes in the two Earth system models (ESMs)
420 that provided input-forcing data to the FishMIP models (Figure S1; Séférian et al., 2020;
421 Tittensor et al., 2021). Across CMIP5 and CMIP6, both ESMs captured observed long-term
422 mean sea-surface temperature (a driver of both models) at the LME scale (Figure S1c, d).
423 Although both models were 0.2-1°C warmer in CMIP6 than CMIP5, they did not show much
424 improvement in resolving net primary production (NPP; a BOATS environmental forcing)
425 at the LME scale (Figure S1a, b). However, across both models, phytoplankton carbon (an
426 EcoOcean forcing) was lower in CMIP6 compared to CMIP5 (Figure S2), which may partly
427 explain why EcoOcean catch bias improved between the two simulations. Ultimately,
428 changes in both MEMs and updated ESM forcings likely contribute to improvements in
429 model skill. Disentangling these drivers would be a fruitful avenue of future research to
430 improve catch and biomass simulations from MEMs.

431 **5.1 Reasons for bias in simulated catches**

432 Bias in fish catch across models are likely driven by a range of factors. For instance, lower
433 trophic level (LTL) biomass and production from ESMs are major drivers of the projected
434 spatial distribution of fish biomass and fisheries catches (Chassot et al., 2010; Heneghan et
435 al., 2021; Kwiatkowski et al., 2020; Laufkötter et al., 2015; Stock et al., 2017; Tagliabue et
436 al., 2021). Thus, discrepancies between observed and modelled LTL variables at the LME
437 scale (Figure S1a, b) will have an impact on MEM fish biomass and therefore catches. In the
438 case of EcoOcean, the one-way forcing of phytoplankton biomass from ESM estimates

439 potentially allows for bias in the estimation of higher trophic level (HTL) biomass,
440 compared to what could be supported by LTL in a two-way coupling. In the case of BOATS,
441 a single energy pathway connects NPP to the accumulation of commercially exploited fish
442 species, while in reality surface pelagic species and bottom living species rely on different
443 food chains, and experience different water temperatures (van Denderen et al., 2018). This
444 may lead to an overestimation of demersal biomass in BOATS, ultimately leading to excess
445 global catches (Guet et al., 2024). Finally, internal climate variability within the ESMs used
446 here was not calibrated to observed variability. This means that seasonal and annual
447 climate patterns affecting simulated catches from the MEMs will not match observation
448 over these shorter time scales.

449 The ecological and fishing components of BOATS and EcoOcean, as with all MEMs, are
450 highly simplified representations of the real world that also differ between models. It is
451 important to note that catch reconstructions are also approximations of reality, subject to
452 numerous sources of bias and error. This necessarily results in discrepancies between
453 models, catch reconstructions and actual catches. For example, fishing in BOATS is here
454 determined by globally homogenous and historically constant economic factors, such as
455 constant fish price and fishing costs, and the development of fisheries in each grid cell is
456 driven by the assumption of a historical increase at a constant rate of the technology-
457 driven catchability of fish biomass, which turns initially unprofitable fishing grounds into
458 profitable ones that can be exploited. This assumption is completed with the assumption of
459 open, unregulated fishing access in the world's oceans (Carozza et al., 2016, 2017). These
460 simplifications can capture broad trends in catches at a global level (Galbraith et al., 2017;
461 Guet et al., 2020), which show a steep increase until the mid-1990s and a later plateau or

462 decline due to overexploitation of fish stocks, and in space a sequential shift of the
463 development of fisheries from cool and productive to warm and unproductive regions
464 (Pauly & Zeller, 2016; Watson & Tidd, 2018). However, they lack important dynamics,
465 which may restrict or change patterns in fishing effort and therefore catches in the real
466 world, particularly at the LME scale considered here.

467 Finally, other internal issues regarding the way key bio-ecological processes are
468 parameterised can lead to divergence between observed and modelled catch rates, as well
469 as the wide differences between BOATS and EcoOcean historical simulations. To identify
470 specific internal drivers of bias and errors, further experimental attribution simulations
471 would be necessary to separate the impact of individual ecological and fisheries
472 components within each model (Steenbeek et al., this issue). Such experimental studies
473 have already successfully identified key drivers of structural uncertainty in the FishMIP
474 ensemble (Heneghan et al., 2021), and biogeochemical modelling community (Laufkötter et
475 al., 2015).

476 **5.2 Evaluation of metrics for marine model assessment**

477 Our case study highlights how the use of multiple metrics is necessary to obtain a multi-
478 faceted perspective of the credibility and reliability of model projections. While summary
479 statistics correlation (R), reliability index (RI) and modelling efficiency (MEF) provide
480 quick and useful information about model credibility, allow comparison between models or
481 regions, and are generally easy to interpret, some important information about the
482 ecosystem is necessarily lost as these metrics reduce the time-series into a single datum
483 (Bennett et al., 2013; Stow et al., 2009). In addition, these metrics do not assess the ability

484 of the models to capture observed geographical patterns in catches, with this shortcoming
485 highlighting the need for spatial skill assessment tools, such as pattern correlation tools, to
486 be developed and applied in parallel.

487 Providing too many metrics that measure the same aspect of model skill may instil false
488 confidence in model assessment by unnecessarily replicating the same result (Olsen et al.,
489 2016). Root mean squared error (RMSE), mean absolute error (MAE) and average error
490 (AE) are all slightly different ways of calculating the same measure – the magnitude of the
491 bias between model simulations and observations. Therefore, it is not necessary to use all
492 three bias metrics moving forward (AE, MAE, RMSE). As AE can be affected by positive and
493 negative discrepancies cancelling each other out, and as RMSE is more sensitive to outliers
494 than MAE, we recommend that MAE be the primary metric used to measure bias. On the
495 other hand, evaluating too few metrics, or metrics that only explore one component of
496 model performance can also instil false confidence. For instance, in previous FishMIP
497 syntheses, model outputs were normalised to show only relative changes (Heneghan et al.,
498 2021; Lotze et al., 2019; Tittensor et al., 2021). Normalisation made it possible to generate
499 a more consistent picture of climate change impacts on marine animal biomass, but may
500 have given a false impression of model agreement since it omitted information on
501 discrepancies in absolute biomass across models. Therefore, we argue that it is essential for
502 metrics exploring both absolute and normalised quantities, as used in this case study, to be
503 deployed when assessing MEM performance.

504 **5.3 Future Research**

505 This paper sets out the features of the CSPS framework and conducts a Level 1 skill
506 assessment for two models within the FishMIP ensemble, on simulated catch. We finish by
507 discussing how this process can continue to be improved.

508 Although 9 global MEMs contribute to the FishMIP ensemble, fisheries catch simulations
509 from only two FishMIP models were available for this assessment. However, we expect that
510 other global and regional models will provide catch outputs as part of the current round of
511 simulations (Blanchard et al., 2024; Frieler et al., 2024) and of ongoing FishMIP efforts to
512 design and implement socioeconomic scenarios that consistently simulate fisheries catch
513 across ecosystem models (Blanchard et al., 2024; Maury et al., this issue).

514 Individual FishMIP models provide a range of integrated outputs besides total catches
515 analysed here, including catches by functional group (i.e., demersal and pelagic) and by size
516 class (e.g., small, medium and large) (Carozza et al., 2016; Cheung et al., 2011; Coll et al.,
517 2020; Maury, 2010; Maury & Poggiale, 2013; Petrik et al., 2019). The analysis of these
518 outputs can add to the model performance picture and can provide insights into modelled
519 ecosystem structure and function. For example, the collapse of large target species and the
520 increase of smaller species due to predation release (Blanchard et al., 2012; Christensen et
521 al., 2014) can drive fisheries catch, but this process is hidden when considering aggregated
522 biomass and catch outputs. Analysis of these existing outputs is an important next step for
523 FishMIP, and forms part of Level 2 (process) and Level 3 (system) assessment in the CSPS
524 framework (Hipsey et al., 2020). Looking ahead, assessing emergent and system-level
525 relationships between ESM variables and MEM output, or between the internal state

526 variables within the MEMs, also offer considerable potential for enhancing Level 2 and
527 Level 3 assessments of MEM performance (Novaglio et al., this issue). Ultimately, an
528 extensive Level 2 and 3 assessment of the FishMIP ensemble will require models to provide
529 outputs that are not currently part of the CMIP and FishMIP protocols, including primary
530 and secondary production rates, biodiversity turnover, trophic transfer rates or growth
531 rates. Eliciting this information in future simulation protocols is therefore critical, since it
532 will provide scope for in-depth assessment of modelled processes across the FishMIP
533 ensemble.

534 The current simulation round of FishMIP is focussed on “Detection, Attribution &
535 Evaluation” (ISIMIP3a, www.fishmip.org), and therefore aims to tackle issues such as
536 resolution, coastal processes, and standardisation of fishing inputs across models. To that
537 end, finer scale inputs from ESMs may help the performance of MEMs at the regional scale.
538 There also exists the opportunity to use FishMIP simulations coupled to ESM models forced
539 by reanalysis data (Blanchard et al., 2024) which are constrained by observational
540 products of atmospheric drivers, to calibrate MEMs, or to conduct post-hoc correction of
541 FishMIP outputs (Gómara et al., 2021; Maury et al., this issue). Unlike fully-coupled ESM
542 historical simulations, ocean-only reanalysis-based simulations would have climate
543 oscillations like ENSO cycles occurring at the correct times in history, and thus would
544 hopefully produce more skillful comparisons of time series (e.g. Barrier et al., 2023).

545 **6 Conclusions**

546 Performing model skill assessment on complex end-to-end ecosystem models is an
547 essential, yet challenging task, and there is still considerable progress to be made before

548 model simulations replicate historical observations. MEMs play an important role in
549 developing our understanding of climate change impacts on future fisheries catches and
550 marine ecosystems, and how that might affect global food security (Blanchard et al., 2012,
551 2017; Booth et al., 2017; Cheung et al., 2010; Cinner et al., 2022; Hollowed et al., 2013).
552 Rigorous ensemble model skill assessment increases confidence in using MEM projections
553 to inform policy, as well as identifying priority areas for future model improvement.

554 Overall, this case study showed that global fishery catch estimates are well correlated with
555 observed trends over time, but both models show important scale mismatches that require
556 further attention. This exercise provides useful information on the performance of two
557 global models contributing to FishMIP and can be further used to drive model development
558 to improve the reliability of climate impact projections, as well as applied more broadly
559 across the whole suite of FishMIP models to enhance the utility of FishMIP as a whole. We
560 finish with a set of summary recommendations for how FishMIP (and other ensemble
561 model projects) could better integrate model ensemble Level 0-3 skill assessment for
562 future simulation protocols:

- 563 1. Level 0: A comprehensive understanding of the underlying assumptions and
564 parameterisations across the model ensemble is essential to understand why MEMs
565 agree or disagree under different conditions. Future protocols targeted at
566 disentangling sources of structural uncertainty across the FishMIP ensemble would
567 concretely improve our understanding of why MEMs behave the way they do. This
568 also includes simulation studies focussed on improving our understanding of the
569 linkages and dependencies between MEMs and the ESMs that force them.

- 570 2. Level 1: FishMIP should move beyond exploring only relative change in simulated
571 variables across the model ensemble, to assessing absolute change and variability.
572 This will require using assessment metrics that capture model bias, such as MAE, RI,
573 or MEF.
- 574 3. Level 2 and 3: To properly assess the processes and emergent properties of
575 ensemble MEMs, future simulation protocols must require modellers to provide
576 more than just aggregate biomass or catch. At the same time, data products on
577 emergent ecosystem properties such as biomass size-spectra need to be assembled
578 at spatial and temporal resolutions appropriate for comparison with global MEMs.

579 The CSPS framework provides a solid basis for standardising skill assessment for FishMIP,
580 to which other metrics (e.g., size-based metrics) could be added. The hierarchical structure
581 and focus of each level act as clear guidelines to measure the predictive validity of MEMs.
582 These initial results show that, although we are yet to fully assess the current ensemble of
583 global marine models, we have the tools and knowledge to tackle this task.

584 **Acknowledgments**

585 MC acknowledges support from the Spanish project ProOceans (Plan Estatal de
586 Investigación Científica, Técnica y de Innovación, 2020, PID2020-118097RB-I00) and the
587 the ‘Severo Ochoa Centre of Excellence’ accreditation (CEX2019-000928-S) to the Institute
588 of Marine Science (ICM-CSIC). CMP acknowledges support from NOAA grants
589 NA200AR4310438, NA200AR4310441, and NA200AR4310442. KOC acknowledges
590 support from the National Research Foundation of South Africa (grant 136481). OM

591 acknowledges support from the European Union's Horizon 2020 research and innovation
592 program under grant agreement N° 817806.

593 **Open Research**

594 All outputs from the FishMIP ensemble are available at
595 www.isimip.org/gettingstarted/data-access/. This includes outputs used for the two
596 models presented in this study, except for EcoOcean outputs from ISIMIP3b (available
597 here: 10.5281/zenodo.11081600).

598 The global fishing catch datasets used for observations are available at
599 <http://dx.doi.org/10.25959/5c522cadbea37> for Watson & Tidd (2018). The Sea Around
600 Us Project dataset is available here: <https://www.seaaroundus.org/data/#/search>).

601 A repository for all R code used to create data visualisations is available on Github here:
602 <https://github.com/nina-rynnne/FishMIP-modelskill>.

603

604 **References**

605 Allen, J. I., & Somerfield, P. J. (2009). A multivariate approach to model skill assessment.
606 *Journal of Marine Systems*, 76(1–2), 83–94.

607 <https://doi.org/10.1016/J.JMARSYS.2008.05.009>

608 Barrier, N., Lengaigne, M., Rault, J., Person, R., Ethé, C., Aumont, O., & Maury, O. (2023).

609 Mechanisms underlying the epipelagic ecosystem response to ENSO in the
610 equatorial Pacific ocean. *Progress in Oceanography*, 213, 103002.

611 <https://doi.org/10.1016/j.pocean.2023.103002>

- 612 Baumberger, C., Knutti, R., & Hirsch Hadorn, G. (2017). Building confidence in climate
613 model projections: An analysis of inferences from fit. *Wiley Interdiscip. Rev. Clim.
614 Change*, 8(3), e454. <https://doi.org/10.1002/wcc.454>
- 615 Bennett, N. D., Croke, B. F. W., Guariso, G., Guillaume, J. H. A., Hamilton, S. H., Jakeman, A. J.,
616 Marsili-Libelli, S., Newham, L. T. H., Norton, J. P., Perrin, C., Pierce, S. A., Robson, B.,
617 Seppelt, R., Voinov, A. A., Fath, B. D., & Andreassian, V. (2013). Characterising
618 performance of environmental models. *Environmental Modelling & Software*, 40, 1–
619 20. <https://doi.org/10.1016/J.ENVSOFT.2012.09.011>
- 620 Blanchard, J. L., Jennings, S., Holmes, R., Harle, J., Merino, G., Allen, J. I., Holt, J., Dulvy, N. K., &
621 Barange, M. (2012). Potential consequences of climate change for primary
622 production and fish production in large marine ecosystems. *Philosophical
623 Transactions of the Royal Society of London. Series B, Biological Sciences*, 367(1605),
624 2979–2989. <https://doi.org/10.1098/RSTB.2012.0231>
- 625 Blanchard, J. L., Novaglio, C., Maury, O., Harrison, C. S., Petrik, C. M., Arcos, L. D. F., Ortega-
626 Cisneros, K., Bryndum-Buchholz, A., Eddy, T., Heneghan, R., Roberts, K. E., Schewe, J.,
627 Bianchi, D., Guiet, J., Denderen, D. V., Palacios-Abrantes, J., Liu, X., Stock, C. A. A.,
628 Rousseau, Y., ... Tittensor, D. (2024). *Detecting, attributing, and projecting global
629 marine ecosystem and fisheries change: FishMIP 2.0.*
630 <https://doi.org/10.22541/essoar.170594183.33534487/v1>
- 631 Blanchard, J. L., Watson, R. A., Fulton, E. A., Cottrell, R. S., Nash, K. L., Bryndum-Buchholz, A.,
632 Büchner, M., Carozza, D. A., Cheung, W. W. L., Elliott, J., Davidson, L. N. K., Dulvy, N. K.,
633 Dunne, J. P., Eddy, T. D., Galbraith, E., Lotze, H. K., Maury, O., Müller, C., Tittensor, D.
634 P., & Jennings, S. (2017). Linked sustainability challenges and trade-offs among

635 fisheries, aquaculture and agriculture. *Nature Ecology and Evolution*, 1(9), 1240–
636 1249. <https://doi.org/10.1038/s41559-017-0258-8>

637 Booth, D. J., Poloczanska, E., Donelson, J. M., Molinos, J. G., & Burrows, M. (2017).
638 Biodiversity and Climate Change in the Oceans. In *Climate Change Impacts on*
639 *Fisheries and Aquaculture* (pp. 63–89). John Wiley & Sons, Ltd.
640 <https://doi.org/10.1002/9781119154051.ch4>

641 Boucher, O., Servonnat, J., Albright, A. L., Aumont, O., Balkanski, Y., Bastrikov, V., Bekki, S.,
642 Bonnet, R., Bony, S., Bopp, L., Braconnot, P., Brockmann, P., Cadule, P., Caubel, A.,
643 Cheruy, F., Codron, F., Cozic, A., Cugnet, D., D'Andrea, F., ... Vuichard, N. (2020).
644 Presentation and Evaluation of the IPSL-CM6A-LR Climate Model. *Journal of*
645 *Advances in Modeling Earth Systems*, 12(7), e2019MS002010.
646 <https://doi.org/10.1029/2019MS002010>

647 Buytaert, W. (2022). *topmodel: Implementation of the Hydrological Model TOPMODEL in R*
648 (0.7.4) [Computer software]. <https://github.com/ICHydro/topmodel>

649 Carozza, D. A., Bianchi, D., & Galbraith, E. D. (2016). The ecological module of BOATS-1.0: A
650 bioenergetically constrained model of marine upper trophic levels suitable for
651 studies of fisheries and ocean biogeochemistry. *Geoscientific Model Development*,
652 9(4), 1545–1565. <https://doi.org/10.5194/GMD-9-1545-2016>

653 Carozza, D. A., Bianchi, D., & Galbraith, E. D. (2017). Formulation, General Features and
654 Global Calibration of a Bioenergetically-Constrained Fishery Model. *PLOS ONE*,
655 12(1), e0169763. <https://doi.org/10.1371/journal.pone.0169763>

- 656 Carslaw, D., & Ropkins, K. (2012). openair—An R package for air quality data analysis.
657 *ENVIRONMENTAL MODELLING & SOFTWARE*, 27–28(0), 52–61. [Computer
658 software] <https://doi.org/10.1016/j.envsoft.2011.09.008>
- 659 Chassot, E., Bonhommeau, S., Dulvy, N. K., Mélin, F., Watson, R., Gascuel, D., & Le Pape, O.
660 (2010). Global marine primary production constrains fisheries catches. *Ecology*
661 *Letters*, 13(4), 495–505. <https://doi.org/10.1111/j.1461-0248.2010.01443.x>
- 662 Cheung, W. W. L., Dunne, J., Sarmiento, J. L., & Pauly, D. (2011). Integrating ecophysiology
663 and plankton dynamics into projected maximum fisheries catch potential under
664 climate change in the Northeast Atlantic. *ICES Journal of Marine Science*, 68(6),
665 1008–1018. <https://doi.org/10.1093/icesjms/fsr012>
- 666 Cheung, W. W. L., Lam, V. W. Y., Sarmiento, J. L., Kearney, K., Watson, R., Zeller, D., & Pauly,
667 D. (2010). Large-scale redistribution of maximum fisheries catch potential in the
668 global ocean under climate change. *Global Change Biology*, 16(1), 24–35.
669 <https://doi.org/10.1111/j.1365-2486.2009.01995.x>
- 670 Christensen, V., Coll, M., Buszowski, J., Cheung, W. W. L., Frölicher, T., Steenbeek, J., Stock, C.
671 A., Watson, R. A., & Walters, C. J. (2015). The global ocean is an ecosystem:
672 Simulating marine life and fisheries. *Global Ecology and Biogeography*, 24(5), 507–
673 517. <https://doi.org/10.1111/geb.12281>
- 674 Christensen, V., Coll, M., Piroddi, C., Steenbeek, J., Buszowski, J., & Pauly, D. (2014). A
675 century of fish biomass decline in the ocean. *Marine Ecology Progress Series*, 512,
676 155–166. <https://doi.org/10.3354/meps10946>

- 677 Christensen, V., & Walters, C. J. (2004). Ecopath with Ecosim: Methods, capabilities and
678 limitations. *Ecological Modelling*, *172*(2), 109–139.
679 <https://doi.org/10.1016/j.ecolmodel.2003.09.003>
- 680 Cinner, J. E., Caldwell, I. R., Thiault, L., Ben, J., Blanchard, J. L., Coll, M., Diedrich, A., Eddy, T.
681 D., Everett, J. D., Folberth, C., Gascuel, D., Guet, J., Gurney, G. G., Heneghan, R. F.,
682 Jägermeyr, J., Jiddawi, N., Lahari, R., Kuange, J., Liu, W., ... Pollnac, R. (2022). Potential
683 impacts of climate change on agriculture and fisheries production in 72 tropical
684 coastal communities. *Nature Communications*, *13*(1), Article 1.
685 <https://doi.org/10.1038/s41467-022-30991-4>
- 686 Coll, M., Steenbeek, J., Pennino, M. G., Buszowski, J., Kaschner, K., Lotze, H. K., Rousseau, Y.,
687 Tittensor, D. P., Walters, C., Watson, R. A., & Christensen, V. (2020). Advancing Global
688 Ecological Modeling Capabilities to Simulate Future Trajectories of Change in
689 Marine Ecosystems. *Frontiers in Marine Science*, *7*, 741.
690 <https://doi.org/10.3389/FMARS.2020.567877/BIBTEX>
- 691 Dunne, J. P., Horowitz, L. W., Adcroft, A. J., Ginoux, P., Held, I. M., John, J. G., Krasting, J. P.,
692 Malyshev, S., Naik, V., Paulot, F., Shevliakova, E., Stock, C. A., Zadeh, N., Balaji, V.,
693 Blanton, C., Dunne, K. A., Dupuis, C., Durachta, J., Dussin, R., ... Zhao, M. (2020). The
694 GFDL Earth System Model Version 4.1 (GFDL-ESM 4.1): Overall Coupled Model
695 Description and Simulation Characteristics. *Journal of Advances in Modeling Earth
696 Systems*, *12*(11), e2019MS002015. <https://doi.org/10.1029/2019MS002015>
- 697 Dunne, J. P., John, J. G., Adcroft, A. J., Griffies, S. M., Hallberg, R. W., Shevliakova, E., Stouffer,
698 R. J., Cooke, W., Dunne, K. A., Harrison, M. J., Krasting, J. P., Malyshev, S. L., Milly, P. C.
699 D., Phillipps, P. J., Sentman, L. T., Samuels, B. L., Spelman, M. J., Winton, M.,

700 Wittenberg, A. T., & Zadeh, N. (2012). GFDL's ESM2 Global Coupled Climate–Carbon
701 Earth System Models. Part I: Physical Formulation and Baseline Simulation
702 Characteristics. *Journal of Climate*, *25*(19), 6646–6665.
703 <https://doi.org/10.1175/JCLI-D-11-00560.1>

704 Eyring, V., Cox, P. M., Flato, G. M., Gleckler, P. J., Abramowitz, G., Caldwell, P., Collins, W. D.,
705 Gier, B. K., Hall, A. D., Hoffman, F. M., Hurtt, G. C., Jahn, A., Jones, C. D., Klein, S. A.,
706 Krasting, J. P., Kwiatkowski, L., Lorenz, R., Maloney, E., Meehl, G. A., ... Williamson, M.
707 S. (2019). Taking climate model evaluation to the next level. *Nature Climate Change*
708 *2019 9:2*, *9*(2), 102–110. <https://doi.org/10.1038/s41558-018-0355-y>

709 FAO. (2022). *The State of World Fisheries and Aquaculture 2022. Towards Blue*
710 *Transformation*. FAO. <https://doi.org/10.4060/cc0461en>

711 Frieler, K., Lange, S., Piontek, F., Reyer, C. P. O., Schewe, J., Warszawski, L., Zhao, F., Chini, L.,
712 Denvil, S., Emanuel, K., Geiger, T., Halladay, K., Hurtt, G., Mengel, M., Murakami, D.,
713 Ostberg, S., Popp, A., Riva, R., Stevanovic, M., ... Yamagata, Y. (2017). Assessing the
714 impacts of 1.5 °C global warming – simulation protocol of the Inter-Sectoral Impact
715 Model Intercomparison Project (ISIMIP2b). *Geoscientific Model Development*, *10*(12),
716 4321–4345. <https://doi.org/10.5194/gmd-10-4321-2017>

717 Frieler, K., Volkholz, J., Lange, S., Schewe, J., Mengel, M., del Rocío Rivas López, M., Otto, C.,
718 Reyer, C. P. O., Karger, D. N., Malle, J. T., Treu, S., Menz, C., Blanchard, J. L., Harrison, C.
719 S., Petrik, C. M., Eddy, T. D., Ortega-Cisneros, K., Novaglio, C., Rousseau, Y., ...
720 Bechtold, M. (2024). Scenario setup and forcing data for impact model evaluation
721 and impact attribution within the third round of the Inter-Sectoral Impact Model

- 722 Intercomparison Project (ISIMIP3a). *Geoscientific Model Development*, 17(1), 1–51.
723 <https://doi.org/10.5194/gmd-17-1-2024>
- 724 Fulton, E. A., Blanchard, J. L., Melbourne-Thomas, J., Plagányi, É. E., & Tulloch, V. J. D. (2019).
725 Where the Ecological Gaps Remain, a Modelers' Perspective. *Frontiers in Ecology and*
726 *Evolution*, 7, 424. <https://doi.org/10.3389/FEVO.2019.00424/BIBTEX>
- 727 Galbraith, E. D., Carozza, D. A., & Bianchi, D. (2017). A coupled human-Earth model
728 perspective on long-term trends in the global marine fishery. *Nature*
729 *Communications* 2017 8:1, 8(1), 1–7. <https://doi.org/10.1038/ncomms14884>
- 730 Geary, W. L., Bode, M., Doherty, T. S., Fulton, E. A., Nimmo, D. G., Tulloch, A. I. T., Tulloch, V. J.
731 D., & Ritchie, E. G. (2020). A guide to ecosystem models and their environmental
732 applications. *Nature Ecology & Evolution* 2020 4:11, 4(11), 1459–1471.
733 <https://doi.org/10.1038/s41559-020-01298-8>
- 734 Gómara, I., Rodríguez-Fonseca, B., Mohino, E., Losada, T., Polo, I., & Coll, M. (2021). Skillful
735 prediction of tropical Pacific fisheries provided by Atlantic Niños. *Environmental*
736 *Research Letters*, 16(5), 054066. <https://doi.org/10.1088/1748-9326/abfa4d>
- 737 Guiet, J., Bianchi, D., Scherrer, K., Heneghan, R., & Galbraith, E. D. (2024). *Small fish biomass*
738 *limits the catch potential in the High Seas*.
739 <https://doi.org/10.22541/au.170967563.32290483/v1>
- 740 Guiet, J., Galbraith, E. D., Bianchi, D., & Cheung, W. W. L. (2020). Bioenergetic influence on
741 the historical development and decline of industrial fisheries. *ICES Journal of Marine*
742 *Science*, 77(5), 1854–1863. <https://doi.org/10.1093/icesjms/fsaa044>
- 743 Halpern, B. S., Walbridge, S., Selkoe, K. A., Kappel, C. V., Micheli, F., D'Agrosa, C., Bruno, J. F.,
744 Casey, K. S., Ebert, C., Fox, H. E., Fujita, R., Heinemann, D., Lenihan, H. S., Madin, E. M.

- 745 P., Perry, M. T., Selig, E. R., Spalding, M., Steneck, R., & Watson, R. (2008). A Global
746 Map of Human Impact on Marine Ecosystems. *Science*, 319(5865), 948–952.
747 <https://doi.org/10.1126/science.1149345>
- 748 Hamilton, S. H., Fu, B., Guillaume, J. H. A., Badham, J., Elsayah, S., Gober, P., Hunt, R. J.,
749 Iwanaga, T., Jakeman, A. J., Ames, D. P., Curtis, A., Hill, M. C., Pierce, S. A., & Zare, F.
750 (2019). A framework for characterising and evaluating the effectiveness of
751 environmental modelling. *Environmental Modelling & Software*, 118, 83–98.
752 <https://doi.org/10.1016/J.ENVSOFT.2019.04.008>
- 753 Hamner, B., Frasco, M., & LeDell, E. (2018). *Metrics: Evaluation Metrics for Machine Learning*
754 (0.1.4) [Computer software]. <https://CRAN.R-project.org/package=Metrics>
- 755 Hatton, I. A., Heneghan, R. F., Bar-On, Y. M., & Galbraith, E. D. (2021). The global ocean size
756 spectrum from bacteria to whales. *Science Advances*, 7(46), 3732.
757 https://doi.org/10.1126/SCIADV.ABH3732/SUPPL_FILE/SCIADV.ABH3732_DATA_
758 [S1.ZIP](https://doi.org/10.1126/SCIADV.ABH3732/SUPPL_FILE/SCIADV.ABH3732_DATA_S1.ZIP)
- 759 Heneghan, R. F., Galbraith, E., Blanchard, J. L., Harrison, C., Barrier, N., Bulman, C., Cheung,
760 W., Coll, M., Eddy, T. D., Erauskin-Extramiana, M., Everett, J. D., Fernandes-Salvador,
761 J. A., Gascuel, D., Guiet, J., Maury, O., Palacios-Abrantes, J., Petrik, C. M., du Pontavice,
762 H., Richardson, A. J., ... Tittensor, D. P. (2021). Disentangling diverse responses to
763 climate change among global marine ecosystem models. *Progress in Oceanography*,
764 198, 102659. <https://doi.org/10.1016/J.POCEAN.2021.102659>
- 765 Hipsey, M. R., Gal, G., Arhonditsis, G. B., Carey, C. C., Elliott, J. A., Frassl, M. A., Janse, J. H., de
766 Mora, L., & Robson, B. J. (2020). A system of metrics for the assessment and

- 767 improvement of aquatic ecosystem models. *Environmental Modelling & Software*,
768 128, 104697. <https://doi.org/10.1016/J.ENVSOFT.2020.104697>
- 769 Hollowed, A. B., Barange, M., Beamish, R. J., Brander, K., Cochrane, K., Drinkwater, K.,
770 Foreman, M. G. G., Hare, J. A., Holt, J., Ito, S. I., Kim, S., King, J. R., Loeng, H., Mackenzie,
771 B. R., Mueter, F. J., Okey, T. A., Peck, M. A., Radchenko, V. I., Rice, J. C., ... Yamanaka, Y.
772 (2013). Projected impacts of climate change on marine fish and fisheries. *ICES*
773 *Journal of Marine Science*, 70(5), 1023–1037.
774 <https://doi.org/10.1093/ICESJMS/FST081>
- 775 Jachner, S., Boogaart, G., & Petzoldt, T. (2007). Statistical Methods for the Qualitative
776 Assessment of Dynamic Models with Time Delay (R Package qualV). *Journal of*
777 *Statistical Software*, 22(8), 1–30.[Computer software]
778 <https://doi.org/10.18637/jss.v022.i08>
- 779 Jakeman, A. J., Letcher, R. A., & Norton, J. P. (2006). Ten iterative steps in development and
780 evaluation of environmental models. *Environmental Modelling & Software*, 21(5),
781 602–614. <https://doi.org/10.1016/J.ENVSOFT.2006.01.004>
- 782 Jennings, S., & Collingridge, K. (2015). Predicting Consumer Biomass, Size-Structure,
783 Production, Catch Potential, Responses to Fishing and Associated Uncertainties in
784 the World's Marine Ecosystems. *PLOS ONE*, 10(7), e0133794.
785 <https://doi.org/10.1371/journal.pone.0133794>
- 786 Kubicek, A., Jopp, F., Breckling, B., Lange, C., & Reuter, H. (2015). Context-oriented model
787 validation of individual-based models in ecology: A hierarchically structured
788 approach to validate qualitative, compositional and quantitative characteristics.

789 *Ecological Complexity*, 22, 178–191.

790 <https://doi.org/10.1016/J.ECOCOM.2015.03.005>

791 Kwiatkowski, L., Torres, O., Bopp, L., Aumont, O., Chamberlain, M., Christian, J. R., Dunne, J.

792 P., Gehlen, M., Ilyina, T., John, J. G., Lenton, A., Li, H., Lovenduski, N. S., Orr, J. C.,

793 Palmieri, J., Santana-Falcón, Y., Schwinger, J., Séférian, R., Stock, C. A., ... Ziehn, T.

794 (2020). Twenty-first century ocean warming, acidification, deoxygenation, and

795 upper-ocean nutrient and primary production decline from CMIP6 model

796 projections. *Biogeosciences*, 17(13), 3439–3470. [https://doi.org/10.5194/bg-17-](https://doi.org/10.5194/bg-17-3439-2020)

797 3439-2020

798 Laufkötter, C., Vogt, M., Gruber, N., Aita-Noguchi, M., Aumont, O., Bopp, L., Buitenhuis, E.,

799 Doney, S. C., Dunne, J., Hashioka, T., Hauck, J., Hirata, T., John, J., Le Quéré, C., Lima, I.

800 D., Nakano, H., Seferian, R., Totterdell, I., Vichi, M., & Völker, C. (2015). Drivers and

801 uncertainties of future global marine primary production in marine ecosystem

802 models. *Biogeosciences*, 12(23), 6955–6984. [https://doi.org/10.5194/bg-12-6955-](https://doi.org/10.5194/bg-12-6955-2015)

803 2015

804 Legates, D. R., & McCabe Jr., G. J. (1999). Evaluating the use of “goodness-of-fit” Measures in

805 hydrologic and hydroclimatic model validation. *Water Resources Research*, 35(1),

806 233–241. <https://doi.org/10.1029/1998WR900018>

807 Lotze, H. K., Tittensor, D. P., Bryndum-Buchholz, A., Eddy, T. D., Cheung, W. W. L., Galbraith,

808 E. D., Barange, M., Barrier, N., Bianchi, D., Blanchard, J. L., Bopp, L., Büchner, M.,

809 Bulman, C. M., Carozza, D. A., Christensen, V., Coll, M., Dunne, J. P., Fulton, E. A.,

810 Jennings, S., ... Worm, B. (2019). Global ensemble projections reveal trophic

811 amplification of ocean biomass declines with climate change. *Proceedings of the*

- 812 *National Academy of Sciences of the United States of America*, 116(26), 12907–12912.
813 https://doi.org/10.1073/PNAS.1900194116/SUPPL_FILE/PNAS.1900194116.SAPP.
814 PDF
- 815 Maureaud, A., Kitchel, Z., Fredston, A., Guralnick, R., Abrantes, J. P., Palomares, D., Pinsky, M.,
816 Shackell, N., Thorson, J., Merigot, B., & Consortium, F. (2023). *FISHGLOB: A*
817 *collaborative infrastructure for marine science and management*. OSF.
818 <https://doi.org/10.31219/osf.io/mh46b>
- 819 Maury, O. (2010). An overview of APECOSM, a spatialized mass balanced “Apex Predators
820 ECOSystem Model” to study physiologically structured tuna population dynamics in
821 their ecosystem. *Progress in Oceanography*, 84(1–2), 113–117.
822 <https://doi.org/10.1016/J.POCEAN.2009.09.013>
- 823 Maury, O., & Poggiale, J.-C. (2013). From individuals to populations to communities: A
824 dynamic energy budget model of marine ecosystem size-spectrum including life
825 history diversity. *Journal of Theoretical Biology*, 324, 52–71.
826 <https://doi.org/10.1016/j.jtbi.2013.01.018>
- 827 Maury, O., Shin, Y.-J., Faugeras, B., Ben Ari, T., & Marsac, F. (2007). Modeling environmental
828 effects on the size-structured energy flow through marine ecosystems. Part 2:
829 Simulations. *Progress in Oceanography*, 74(4), 500–514.
830 <https://doi.org/10.1016/j.pocean.2007.05.001>
- 831 Maury et al. (this issue). The Ocean System Pathways: A new scenario framework to
832 investigate the future of ocean ecosystems and fisheries. In Submission.
- 833 Mayer, D. G., & Butler, D. G. (1993). Statistical validation. *Ecological Modelling*, 68(1–2), 21–
834 32. [https://doi.org/10.1016/0304-3800\(93\)90105-2](https://doi.org/10.1016/0304-3800(93)90105-2)

- 835 Novaglio, C., Blanchard, J. L., Plank, M. J., van Putten, E. I., Audzijonyte, A., Porobic, J., &
836 Fulton, E. A. (2022). Exploring trade-offs in mixed fisheries by integrating fleet
837 dynamics into multispecies size-spectrum models. *Journal of Applied Ecology*, 59(3),
838 715–728. <https://doi.org/10.1111/1365-2664.14086>
- 839 Novaglio, C., Bryndum-Buchholz, A., Tittensor, D., Eddy, T. D., Lotze, H. K., Harrison, C. S.,
840 Heneghan, R., Maury, O., Ortega-Cisneros, K., Petrik, C. M., Roberts, K. E., &
841 Blanchard, J. L. (2024). *The Past and Future of the Fisheries and Marine Ecosystem*
842 *Model Intercomparison Project*. ESS Open Archive.
843 <https://doi.org/10.22541/essoar.170542252.20348236/v1>
- 844 Novaglio et al. (this issue). Using ecological theory to evaluate and reduce uncertainties of a
845 marine ecosystem model ensemble. In submission.
- 846 Olsen, E., Fay, G., Gaichas, S., Gamble, R., Lucey, S., & Link, J. S. (2016). Ecosystem Model Skill
847 Assessment. Yes We Can! *PLOS ONE*, 11(1), e0146467.
848 <https://doi.org/10.1371/JOURNAL.PONE.0146467>
- 849 Ortega-Cisneros, K., Cochrane, K., & Fulton, E. (2017). An Atlantis model of the southern
850 Benguela upwelling system: Validation, sensitivity analysis and insights into
851 ecosystem functioning. *ECOLOGICAL MODELLING*, 355, 49–63.
852 <https://doi.org/10.1016/j.ecolmodel.2017.04.009>
- 853 Ortega-Cisneros K., Fierros Arcos D., Novaglio C., Eddy T., Coll M., Fulton E., Oliveros-Ramos
854 R. *et al.* Submitted. An integrated workflow for the detection and attribution of marine
855 ecosystem change, from global to regional scales. *Earth's Future*. This Issue
856 MS#2024EF004826
857

- 858 Pauly, D., & Zeller, D. (2016). Catch reconstructions reveal that global marine fisheries
859 catches are higher than reported and declining. *Nature Communications*, 7.
860 <https://doi.org/10.1038/NCOMMS10244>
- 861 Petrik, C. M., Luo, J. Y., Heneghan, R. F., Everett, J. D., Harrison, C. S., & Richardson, A. J.
862 (2022). Assessment and Constraint of Mesozooplankton in CMIP6 Earth System
863 Models. *Global Biogeochemical Cycles*, 36(11), e2022GB007367.
864 <https://doi.org/10.1029/2022GB007367>
- 865 Petrik, C. M., Stock, C. A., Andersen, K. H., van Denderen, P. D., & Watson, J. R. (2019).
866 Bottom-up drivers of global patterns of demersal, forage, and pelagic fishes.
867 *Progress in Oceanography*, 176, 102124.
868 <https://doi.org/10.1016/j.pocean.2019.102124>
- 869 Planque, B., Aarflot, J. M., Buttay, L., Carroll, J., Fransner, F., Hansen, C., Husson, B.,
870 Langangen, Ø., Lindstrøm, U., Pedersen, T., Primicerio, R., Sivel, E., Skogen, M. D.,
871 Strombom, E., Stige, L. C., Varpe, Ø., & Yoccoz, N. G. (2022). A standard protocol for
872 describing the evaluation of ecological models. *Ecological Modelling*, 471, 110059.
873 <https://doi.org/10.1016/j.ecolmodel.2022.110059>
- 874 Power, M. (1993). The predictive validation of ecological and environmental models.
875 *Ecological Modelling*, 68(1–2), 33–50. [https://doi.org/10.1016/0304-](https://doi.org/10.1016/0304-3800(93)90106-3)
876 [3800\(93\)90106-3](https://doi.org/10.1016/0304-3800(93)90106-3)
- 877 Refsgaard, J. C., Madsen, H., Andréassian, V., Arnbjerg-Nielsen, K., Davidson, T. A., Drews, M.,
878 Hamilton, D. P., Jeppesen, E., Kjellström, E., Olesen, J. E., Sonnenborg, T. O., Trolle, D.,
879 Willems, P., & Christensen, J. H. (2014). A framework for testing the ability of models

880 to project climate change and its impacts. *Climatic Change*, 122(1–2), 271–282.
881 <https://doi.org/10.1007/S10584-013-0990-2/FIGURES/3>

882 Ricard, D., Minto, C., Jensen, O. P., & Baum, J. K. (2012). Examining the knowledge base and
883 status of commercially exploited marine species with the RAM Legacy Stock
884 Assessment Database. *Fish and Fisheries*, 13(4), 380–398.
885 <https://doi.org/10.1111/j.1467-2979.2011.00435.x>

886 Rykiel, E. J. (1996). Testing ecological models: The meaning of validation. *Ecological*
887 *Modelling*, 90(3), 229–244. [https://doi.org/10.1016/0304-3800\(95\)00152-2](https://doi.org/10.1016/0304-3800(95)00152-2)

888 Scherrer, K., & Galbraith, E. (2020). Regulation strength and technology creep play key
889 roles in global long-term projections of wild capture fisheries. *ICES Journal of Marine*
890 *Science*, 77(7–8), 2518–2528. <https://doi.org/10.1093/icesjms/fsaa109>

891 Séférian, R., Berthet, S., Yool, A., Palmiéri, J., Bopp, L., Tagliabue, A., Kwiatkowski, L.,
892 Aumont, O., Christian, J., Dunne, J., Gehlen, M., Ilyina, T., John, J. G., Li, H., Long, M. C.,
893 Luo, J. Y., Nakano, H., Romanou, A., Schwinger, J., ... Yamamoto, A. (2020). Tracking
894 Improvement in Simulated Marine Biogeochemistry Between CMIP5 and CMIP6.
895 *Current Climate Change Reports*, 6(3), 95–119. [https://doi.org/10.1007/s40641-](https://doi.org/10.1007/s40641-020-00160-0)
896 [020-00160-0](https://doi.org/10.1007/s40641-020-00160-0)

897 Sepulchre, P., Caubel, A., Ladant, J.-B., Bopp, L., Boucher, O., Braconnot, P., Brockmann, P.,
898 Cozic, A., Donnadieu, Y., Dufresne, J.-L., Estella-Perez, V., Ethé, C., Fluteau, F., Foujols,
899 M.-A., Gastineau, G., Ghattas, J., Hauglustaine, D., Hourdin, F., Kageyama, M., ... Tardif,
900 D. (2020). IPSL-CM5A2 – an Earth system model designed for multi-millennial
901 climate simulations. *Geoscientific Model Development*, 13(7), 3011–3053.
902 <https://doi.org/10.5194/gmd-13-3011-2020>

- 903 Steenbeek, J., Buszowski, J., Chagaris, D., Christensen, V., Coll, M., Fulton, E. A., Katsanevakis,
904 S., Lewis, K. A., Mazaris, A. D., Macias, D., de Mutsert, K., Oldford, G., Pennino, M. G.,
905 Piroddi, C., Romagnoni, G., Serpetti, N., Shin, Y. J., Spence, M. A., & Stelzenmüller, V.
906 (2021). Making spatial-temporal marine ecosystem modelling better – A
907 perspective. *Environmental Modelling & Software*, *145*, 105209.
908 <https://doi.org/10.1016/J.ENVSOFT.2021.105209>
- 909 Steenbeek et al. (this issue). Making ecosystem modelling operational - a novel distributed
910 execution framework to systematically explore ecological responses to divergent climate
911 trajectories. Paper #2023EF004295.
- 912 Stock, C. A., John, J. G., Rykaczewski, R. R., Asch, R. G., Cheung, W. W. L., Dunne, J. P.,
913 Friedland, K. D., Lam, V. W. Y., Sarmiento, J. L., & Watson, R. A. (2017). Reconciling
914 fisheries catch and ocean productivity. *Proceedings of the National Academy of
915 Sciences of the United States of America*, *114*(8), E1441–E1449.
916 <https://doi.org/10.1073/PNAS.1610238114>
- 917 Stow, C. A., Jolliff, J., McGillicuddy, D. J., Doney, S. C., Allen, J. I., Friedrichs, M. A. M., Rose, K.
918 A., & Wallhead, P. (2009). Skill assessment for coupled biological/physical models of
919 marine systems. *Journal of Marine Systems*, *76*(1–2), 4–15.
920 <https://doi.org/10.1016/J.JMARSYS.2008.03.011>
- 921 Sturludottir, E., Desjardins, C., Elvarsson, B., Fulton, E. A., Gorton, R., Logemann, K., &
922 Stefansson, G. (2018). End-to-end model of Icelandic waters using the Atlantis
923 framework: Exploring system dynamics and model reliability. *Fisheries Research*,
924 *207*, 9–24. <https://doi.org/10.1016/J.FISHRES.2018.05.026>

- 925 Tagliabue, A., Kwiatkowski, L., Bopp, L., Butenschön, M., Cheung, W., Lengaigne, M., &
926 Vialard, J. (2021). Persistent Uncertainties in Ocean Net Primary Production Climate
927 Change Projections at Regional Scales Raise Challenges for Assessing Impacts on
928 Ecosystem Services. *Frontiers in Climate*, 3.
929 <https://www.frontiersin.org/articles/10.3389/fclim.2021.738224>
- 930 Taylor, K. E. (2001). Summarizing multiple aspects of model performance in a single
931 diagram. *Journal of Geophysical Research: Atmospheres*, 106(D7), 7183–7192.
932 <https://doi.org/10.1029/2000JD900719>
- 933 Tittensor, D. P., Eddy, T. D., Lotze, H. K., Galbraith, E. D., Cheung, W., Barange, M., Blanchard,
934 J. L., Bopp, L., Bryndum-Buchholz, A., Büchner, M., Bulman, C., Carozza, D. A.,
935 Christensen, V., Coll, M., Dunne, J. P., Fernandes, J. A., Fulton, E. A., Hobday, A. J.,
936 Huber, V., ... Walker, N. D. (2018). A protocol for the intercomparison of marine
937 fishery and ecosystem models: Fish-MIP v1.0. *Geoscientific Model Development*,
938 11(4), 1421–1442. <https://doi.org/10.5194/GMD-11-1421-2018>
- 939 Tittensor, D. P., Novaglio, C., Harrison, C. S., Heneghan, R. F., Barrier, N., Bianchi, D., Bopp, L.,
940 Bryndum-Buchholz, A., Britten, G. L., Büchner, M., Cheung, W. W. L., Christensen, V.,
941 Coll, M., Dunne, J. P., Eddy, T. D., Everett, J. D., Fernandes-Salvador, J. A., Fulton, E. A.,
942 Galbraith, E. D., ... Blanchard, J. L. (2021). Next-generation ensemble projections
943 reveal higher climate risks for marine ecosystems. *Nature Climate Change* 2021
944 11:11, 11(11), 973–981. <https://doi.org/10.1038/s41558-021-01173-9>
- 945 van Denderen, P. D., Lindegren, M., MacKenzie, B. R., Watson, R. A., & Andersen, K. H.
946 (2018). Global patterns in marine predatory fish. *Nature Ecology & Evolution*, 2(1),
947 Article 1. <https://doi.org/10.1038/s41559-017-0388-z>

- 948 Wagener, T., Reinecke, R., & Pianosi, F. (2022). On the evaluation of climate change impact
949 models. *WIREs Climate Change*, 13(3), e772. <https://doi.org/10.1002/wcc.772>
- 950 Watson, R. (2019). *Global Fisheries Landings V4.0* (4.0) [dataset]. Institute for Marine and
951 Antarctic Studies (IMAS), University of Tasmania (UTAS).
952 <https://doi.org/10.25959/5C522CADBEA37>
- 953 Watson, R., & Tidd, A. (2018). Mapping nearly a century and a half of global marine fishing:
954 1869–2015. *Marine Policy*, 93, 171–177.
955 <https://doi.org/10.1016/j.marpol.2018.04.023>
- 956