

Development and testing of effort-tailored superensemble and two-stage catch-only models for estimation of stock status

Olaf P. Jensen¹ and Christopher M. Free¹

¹ Department of Marine & Coastal Sciences, Rutgers University

Contact: olaf.p.jensen@gmail.com

1. Summary

Catch-only models (COMs) represent a useful class of tools for categorizing the status of the vast majority of the world's fisheries which lack formal stock assessments. Testing of COMs on simulated data highlighted substantial differences in performance depending on the fishing effort (and thus fishing mortality) trend used to generate the simulated data. Because of this, no single method is optimal for all fishing effort trend scenarios. Smaller differences in model performance are also apparent when comparing different simulated life histories.

We developed and tested two types of higher order COMs (i.e., COMs composed of two or more individual COMs): 1. Effort-tailored superensembles, i.e., superensembles fit to simulated catch data from individual effort dynamics scenarios and 2. A mechanistic two-stage COM that used the Zhou-BRT method to constrain an influential parameter of the cMSY model.

The effort-tailored superensembles were consistently the best (or equally good) performers in both continuous and categorical cross-validation using simulated data and data from real assessed stocks. However, even these best models showed only fair performance when applied to data from real stocks with some evidence of a pessimistic bias overall. Of perhaps greater concern, these effort-tailored superensembles perform poorly when the effort trend is mis-specified. The new two-stage COM developed here performed better than existing two-stage COMs, but generally not as well as the superensembles.

In summary, if a COM is to be used for status estimation, and effort trends can be reliably categorized, we recommend the use of Super-4. However, testing against RAMLDB stocks indicates that accuracy is modest and some bias may exist (i.e., a tendency to slightly underestimate B/B_{MSY}). We cannot yet recommend any COMs as a reliable tool for providing management advice without thorough testing in a management strategy evaluation (MSE) framework.

2. Introduction

Catch-only models (COMs) are a stock assessment method that requires no time series data other than catch. Despite ongoing debate about the information content present in a catch time series (Pauly et al. 2013), development, testing, and application of new COMs has preceded rapidly over the past few years (Rosenberg et al. 2014, 2017; Zhou et al. 2016, 2017; Free et al. 2017, Anderson et al. 2017). If these models are determined to be reliable, COMs represent a useful class of tools for categorizing the status of the vast majority of the world's

fisheries which lack formal stock assessments. One especially important application would be to use COMs to improve the accuracy of global stock status classifications in the FAO's State of World Fisheries and Aquaculture (SOFIA) reports.

Testing of COMs by Rosenberg et al. (2014) and Jensen and Free (2017) on simulated catch time series has revealed that model performance varies greatly depending on the simulated effort dynamics. This is not surprising as effort is a “hidden” (i.e., unobserved) variable in all COMs that are based on an underlying population dynamics model (Thorson et al. 2013). Methods such as cMSY, which perform well under some scenarios of simulated change in fishing effort, perform quite poorly under others. No single method is optimal for all scenarios. Thus, there is a strong potential for significant improvements in performance of COMs, if they are tailored to specific effort scenarios and combined with an expert classification of effort dynamics, which could be accomplished through a simple questionnaire given to data providers.

Superensembles, a type of COM which uses the predictions of other COMs as input for a statistical model, have been shown to generally outperform individual COMs at predicting stock status and trends (Anderson et al. 2017). However, the superensembles developed by Anderson et al. (2017) were fit to simulation data aggregated across Rosenberg et al.'s (2014) four different effort dynamics scenarios. This provides a certain degree of robustness to uncertainty in the effort dynamics, but it is sub-optimal for stocks for which effort dynamics are known. The Rosenberg et al. (2014) effort dynamics scenarios are quite general (increasing, decreasing, or stable trend and biomass-coupled effort), and it is likely that someone familiar with a fishery would be able to accurately categorize it in terms of its effort dynamics. Development and testing of superensemble models specific to each of Rosenberg et al.'s (2014) four effort dynamics scenarios will provide a useful set of tools for estimating status for stocks where only a time series of catch and an effort classification are available.

Mechanistic combinations of models also hold promise. For example, performance of cMSY depends on assumptions regarding depletion in the final year, a quantity that is estimated by the Zhou-BRT COM. Thus, a two-stage approach which combines these two methods may perform better than either method alone.

In addition, life history characteristics of the species in question also play a role in determining the accuracy of stock status estimates from COMs and the best performing COM varies somewhat depending on life history – though this difference is not as great as the differences attributable to effort dynamics (Jensen and Free 2017, Table 8).

Our goals were: (1) to fit individual effort-tailored superensemble models to simulated catch data from each of the four effort dynamics scenarios used by Rosenberg et al. (2014); (2) to develop a two-stage COM using Zhou-BRT to estimate the final year depletion used by cMSY; (3) to test these two approaches through cross-validation using withheld simulated data and catch data from real stocks with full assessments.

3. Methods

3.1 Datasets

We developed the effort-tailored superensemble models using simulated fish stocks from Rosenberg et al. (2014) and tested the models on a set of simulated stocks withheld from model training and on real fish stocks in the RAM Legacy Stock Assessment Database (RAMLDB v. 2.95; Ricard et al. 2012). We tested the two-stage catch-only model on both the simulated and RAMLDB stocks. The Rosenberg et al. (2014) simulated stocks represent a fully factorial set of 5760 simulated fisheries comprised of three fish life histories, three levels of initial biomass depletion, four exploitation scenarios, two levels of recruitment variability, two levels of recruitment autocorrelation, and two levels of measurement error, with each combination of parameters run through ten stochastic iterations (**Supp. Table 1**). The RAMLDB is a global database of catch data and stock assessment output, including reference points and time series of biomass and fishing mortality. We used 161 of the 193 RAMLDB stocks used in Free et al. (2017) to allow for the comparison of model performance against the refined ORCS approach (rORCS), a COM that has shown strong performance in estimating stock status. Invertebrate stocks and stocks with catch time series < 20 years long after trimming years of zero catch from the beginning of the time series were excluded for this analysis.

3.2 Superensemble models

3.2.1 Building the superensemble models

We developed three sets of superensemble models to estimate stock status: (1) one model for all stocks (Super-1); (2) four models for stocks experiencing each of four effort dynamics (ED) scenarios (**Supp. Table 2**; Super-4); and (3) twelve models for stocks characterized by each of twelve combinations of three life history (LH) archetypes and four effort dynamics scenarios (**Supp. Table 3**; Super-12).

Each superensemble model uses boosted regression trees (BRT) to estimate stock status (B/B_{MSY}) from the B/B_{MSY} estimates of five individual catch-only assessment models (**Table 1**) and two spectral properties of the catch time series. Boosted regression trees combine regression and machine learning, offer predictive power superior to other modeling methods (Elith et al. 2008), and produced the best superensemble model in Anderson et al. (2017). We included the 0.05 and 0.20 spectral densities of the scaled catch time series (catch divided by maximum catch) because they were shown to improve predictive performance in Anderson et al. (2017). Because B/B_{MSY} is a ratio bounded at zero, we fit the BRT models using the log of B/B_{MSY} and exponentiated predictions from the model. Thus, each of the superensemble models has the following conceptual structure:

$$\log \theta = f(\beta_{CMSY-17}, \beta_{COMSIR}, \beta_{SSCOM}, \beta_{mPRM}, \beta_{zBRT}, SD_{0.05}, SD_{0.20})$$

where θ represents the superensemble estimate of B/B_{MSY} , β 's represent the individual model estimates of B/B_{MSY} , and SD 's are the spectral densities of the scaled catch time series.

We divided the simulated stocks for model training (90% of data) and testing (10% of data) by withholding the 10th iteration of each simulation scenario. The training stocks were used to fit the BRT models while the test stocks were used to independently evaluate each model's predictive ability. For each model, we performed an initial grid search for the BRT model parameters that minimize the RMSE using 10-fold cross validation to avoid overfitting (**Appendix A**). The grid search looked for the optimal number of trees, interaction depth, and learning rate, and the optimal parameters are listed in **Table 2**. The BRT models were fit using the *caret* (Kuhn 2016) and *gbm* (Ridgeway 2016) packages in R v.3.4.2 (R Core Team 2017).

3.2.2 Testing the superensemble models

We evaluated the ability of the BRT models to predict continuous (i.e., B/B_{MSY}) and categorical (i.e., under, fully, or overexploited) stock status through testing on the simulated stocks withheld from model training and on the fully independent RAMLDB stocks. We followed the performance evaluation framework proposed by Jensen and Free (2017), briefly: continuous performance was measured in terms of rank-order correlation, accuracy (median absolute proportional error), and bias (median proportional error) and categorical performance was measured in terms of percent accuracy and Cohen's kappa, which accounts for the probability of correct classification by chance. We compared the performance of the BRT models to 10 of the 11 COMs evaluated by Jensen and Free (2017) (**Table 1**). The original GBM-Superensemble model (Anderson et al. 2017) was excluded because it was trained to predict the mean B/B_{MSY} over the last five years of the catch time series rather than B/B_{MSY} in the final year. Three of these methods – SSP-2002, SSP-2013, and rORCS – predict only categorical stock status.

Because the ED- and LH-ED-tailored BRT models require an expert to specify the underlying effort dynamics, we also evaluated the sensitivity of these models to the misspecification of effort history by comparing performance when misspecified with performance when correctly specified. The sensitivity analysis was only performed on the test simulated stocks and not on the RAMLDB stocks given the limited and unequal representation of the ED scenarios in the RAMLDB stocks (**Supp. Table 3**).

To test the ED- and LH-ED-tailored BRT models on the RAMLDB stocks, we classified the RAMLDB stocks into life history and effort dynamics categories consistent with the simulated stocks (**Supp. Tables 2-3**). We classified life history categories based on taxonomic family according to **Supp. Table 4**. We only classified the RAMLDB stocks as experiencing “constant”, “increasing”, or “roller coaster” dynamics given that “biomass-coupled” effort can yield all three patterns. Classification was automated and proceeded as follows: First, we fit linear and two-slope segmented regression models to the effort time series of the RAMLDB stocks and identified the best regression model for each stock using AIC. We classified stocks whose effort histories were best described by segmented regressions with positive initial slopes and negative final slopes as experiencing “roller coaster” effort (**Appendix B**). The remaining stocks were classified based on the significance and direction of the slope of the linear regression on effort history: stocks with non-significant slopes were classified as experiencing “constant” effort and

stocks with significantly positive or negative slopes were classified as experiencing “increasing” or “decreasing” effort, respectively (**Appendix B**). RAMLDB stocks experiencing “decreasing” effort were reclassified as experiencing “roller coaster” effort based on the assumption that high effort at the beginning of the available time series must have been preceded by an unmodelled period of increasing effort during the fishery development phase (Csirke & Sharp 1984). Overall, we identified 11 constant, 41 increasing, and 109 roller coaster RAMLDB stocks (**Supp. Table 3**).

3.3 Two-stage catch only model

We developed a new two-stage catch only model that uses the Zhou-BRT method (Zhou et al 2017) to inform a final year saturation prior for use in the cMSY-17 stock reduction analysis (Froese et al. 2017; saturation = 1 – depletion = $B / B_0 = B/B_{MSY} / 2$). Currently, cMSY-17 derives a final year saturation prior using simple rules based on the ratio of catch in the final year to the maximum catch where lower and higher catch ratios imply lower and higher saturation, respectively (**Figure 13**). Unfortunately, these priors exhibit low overlap with observed saturations from both the simulated and RAMLDB stocks (**Figure 14**), which is unsurprising given that catch ratios are poor indicators of stock status (Branch et al. 2011; Carruthers et al. 2012; Jensen & Free 2017). For this reason, we use the saturation predictions of the Zhou-BRT method, which are more highly correlated with stock status (Jensen & Free 2017), to inform our final year saturation priors. Specifically, we estimate the mean saturation using the Zhou-BRT method, then identify the lower and upper bounds of a uniform prior using the 95% confidence interval of the skewed normal distributions described in Zhou et al. (2017):

$$\begin{aligned} \text{If } S \leq 0.5: & \quad f(\xi = \max(S_{ZBRT}, 0) - 0.072, \omega = 0.189, \alpha = 0.763) \\ \text{If } S > 0.5: & \quad f(\xi = \max(S_{ZBRT}, 0) + 0.179, \omega = 0.223, \alpha = 0.904) \end{aligned}$$

where ξ is the location parameter, ω is the scale parameter, and α is the skewness parameter and the lower bound is bounded at zero (**Figure 14**). These values are passed directly into the expert-defined lower and upper bounds for the uniform final year saturation in cMSY-17. The resilience values required by cMSY-17 were specified for the simulated and RAMLDB using the methods of Jensen & Free (2017). The performance of this new two-stage catch-only model was evaluated in the same framework described for the superensemble models.

4. Results

4.1 Superensemble model results

The superensemble models were the best overall predictors of both continuous and categorical stock status when evaluated on the test simulated stocks (**Figure 1**) with the effort-tailored superensemble models (Super-4, Super-12) performing better than the overall superensemble model (Super-1; **Figures 1-2**). The LH-ED-tailored models (Super-12), which performed only slightly better than the ED-tailored models (Super-4), exhibited high rank-order correlation, low

inaccuracy, and low bias in predicting B/B_{MSY} and were both “good” and highly accurate classifiers of categorical stock status (**Figure 1**). Although the overall performance of the LH-ED tailored models on the RAMLDB stocks was not as outstanding as on the test simulated stocks, they were still consistently among the best performing methods. They exhibited the highest correlation and third lowest inaccuracy in their predictions of B/B_{MSY} and were the second-best classifiers of stock status after the rORCS approach, which was trained on the RAMLDB (**Figure 1**).

Across the four ED scenarios in the test simulated stocks (i.e., grouping life histories), the LH-ED-tailored models offered the best predictions of B/B_{MSY} and were consistently among the best classifiers of stock status (**Figure 3**). When tested on the RAMLDB stocks, the LH-ED-tailored models offered predictions of B/B_{MSY} that were often the most highly correlated but in the middle of the evaluated COMs in terms of continuous inaccuracy and categorical classification performance (**Figure 4**).

Across LH-ED scenarios in the test simulated stocks, the LH-ED-tailored models exhibited consistently low inaccuracy in B/B_{MSY} predictions and were almost always the least inaccurate of the COMs evaluated. However, they were considerably more mixed in their rank-order correlation. While they were the best predictors for stocks in the LP-Coupled, LP-Increasing, and DE-Coupled scenarios, they were among the worst predictors in the SP-Constant and DE-Roller coaster scenarios (**Figure 5**). The categorical performance of the LH-ED-tailored models across LH-ED scenarios mirrored their continuous performance across scenarios, mainly: (1) they exhibited consistently high accuracy and were almost always the most accurate of the COMs evaluated but were frequently “poor” classifiers of stock status and (2) they performed well in the SP/LP/DE-Coupled and LP-Increasing scenarios (**Figure 6**). The LH-ED-tailored models exhibited a mixed performance when tested on the RAMLDB stocks. In only two scenarios did they offer the most highly correlated predictions of B/B_{MSY} and in nearly all scenarios they exhibited inaccuracy values in the middle of the evaluated COMs (**Figure 7**). Similarly, they generally offered classification performance in the middle of the evaluated COMs (**Figure 8**).

The B/B_{MSY} predictions of both the ED- and LH-ED-tailored superensembles applied to the test simulated stocks become more inaccurate when effort dynamics are misspecified, even if they occasionally become more correlated (**Figures 9-10**). Similarly, the categorical status predictions of both sets of superensembles become less accurate when effort dynamics are misspecified, even if misspecifications occasionally result in a higher Cohen’s kappa (**Figures 9 & 11**). When using the ED-tailored models, there are particularly extreme costs to prediction accuracy when misspecifying stocks experiencing “biomass-coupled” effort and stocks experiencing “increasing” effort (**Figure 9**). When using the LH-ED-tailored models, there are particularly extreme costs to prediction accuracy to misspecifying stocks experiencing “increasing” effort and to misspecifying stocks experiencing “roller coaster” effort as experiencing “constant” effort (**Figures 10-11**). Stocks experiencing “constant” effort are the least sensitive to misspecification when using both sets of models (**Figures 9-11**).

4.2 Two-stage model results

The two-stage catch-only model (2-Stage) performed worse on the test simulated stocks and RAMLDB stocks than the superensemble models (**Figures 1 & 3**) but performed better than the other two-stage models: OCOM and cMSY-17 (**Figures 3 & 14**). OCOM, which uses the same final year saturation prior as the two-stage model but employs a different stock reduction analysis, produces a heavily bimodal distribution of status estimates (**Figure 14**). cMSY-17, which uses the same stock reduction analysis as the two-stage model but with a different final year saturation prior, produces heavily biased status estimates (**Figure 14**). The two-stage model had no standout performances when tested on the ED or LH-ED scenarios in either the simulated or RAMLDB stocks (**Figures 3-8**).

5. Discussion

Although the effort-tailored superensembles are in many ways an improvement over other COMs, including the general superensemble (Super-1), they offer clear benefits only under a constrained set of stock characteristics, and the cross-validation performance of all COMs when applied to data from real stocks suggests reasons for caution. In cross-validation testing against withheld simulated data (Fig. 1), both the effort-tailored (Super-4) and the effort and life history-tailored (Super-12) models perform substantially better than the next best performing model (Super-1). The overall performance differences between Super-4 and Super-12 are minor, suggesting that additional fitting of superensembles to specific life histories offers little advantage over tailoring superensembles to effort dynamics alone. Testing COMs against real stocks from the RAMLDB offers an arguably more realistic assessment of their performance when applied to real data-poor stocks. Here none of the COMs distinguished themselves as consistently reliable and performance was generally similar among a set of several of the best COMs. For categorical performance (Fig. 1, bottom right), Super-4 and Super-12 were the best models. Both had kappa approximately 0.2 (fair) and classification accuracy of 50% (not equivalent to a coin flip as there were three categories: underexploited, fully exploited, and overexploited). While rORCS performed substantially better than other COMs, this represents a different type of comparison since rORCS was fit to (a different set of) stocks from the RAMLDB. For continuous performance (Fig. 1, bottom left), there was no single best model for both metrics, but Super-4 and Super-12 both did comparatively well by both metrics. However, like most of the COMs, they had a negative bias: that is, they tended to slightly underestimate B/B_{MSY} .

Performance of the effort-tailored superensembles differed substantially depending on the underlying effort dynamics of the stock. Several COMs performed well (rank order correlation > 0.5) in continuous prediction of RAMLDB stocks with no significant F trend (the constant effort scenario), with effort-tailored models Super-4 and Super-12 showing the highest rank order correlation (Fig. 4). This is not surprising since if effort is stationary (as well as catchability and productivity), any trend in catch is likely to reflect a trend in biomass. In contrast, although Super-4 and Super-12 showed the highest rank order correlation, none of the COMs had a rank order correlation > 0.5 for either the increasing or roller coaster F stocks in the RAMLDB.

Applying an effort-tailored superensemble to catch data simulated using a different effort scenario always results in poorer performance, though the extent of this performance degradation depends on the specific effort scenario. At the extreme, the superensemble fit to the constant effort scenario performs relatively poorly even on simulated data derived from a constant effort pattern (Fig. 9). When applied to catch data simulated from other effort scenarios both continuous and categorical performance is very poor. Thus, this particular superensemble should only be used when there is high certainty that effort has been constant, and even then predictions are highly uncertain. The best performing of the effort-tailored superensembles is the one fit to simulated data from the biomass coupled effort dynamics scenario. It performs well in both continuous (rank order correlation = 0.81, inaccuracy = 0.11) and categorical (accuracy = 88%, kappa = 0.71) cross-validation (Fig. 9). However, even this model shows poor predictive performance when applied to data generated from other effort scenarios.

The two-stage model that we developed and tested here performs better than other two-stage or individual COMs, though not as well as the superensembles. Two of the other methods, cMSY and OCOM, are also inherently two-stage approaches in which either a model (OCOM) or a simple catch ratio (cMSY) are used to place a constraint (an informative uniform prior) on depletion (or saturation) in the final year. OCOM predictions were bimodal (Fig. 14) with a notable gap in predicted B/B_{MSY} between approximately 1.0 and 1.3. Structurally, OCOM and the two-stage model that we developed are quite similar. Both use Zhou-BRT to constrain depletion (or saturation) in the final year coupled with a stock reduction analysis. The relatively poor performance of all three of these two-stage COMs compared to the superensembles is most likely a result of the strong influence of the depletion constraint combined with an inability to adequately estimate this constraint (Fig. 13).

In summary, if a COM is to be used for status estimation, and effort trends can be reliably categorized, we recommend the use of Super-4. However, testing against RAMLDB stocks indicates that accuracy is modest and some bias may exist (i.e., a tendency to slightly underestimate B/B_{MSY}). We cannot yet recommend any COMs as a reliable tool for providing management advice without thorough testing in a management strategy evaluation (MSE) framework.

References

- Anderson, S.C., J. Afflerbach, A.B. Cooper, M. Dickey-Collas, O.P. Jensen, K.M. Kleisner, C. Longo, G.C. Osio, D. Ovando, C. Minte-Vera, C. Minto, I. Mosqueira, A.A. Rosenberg, E.R. Selig, J.T. Thorson, and J.C. Walsh. 2016. datalimited: Stock assessment methods for data-limited fisheries. R package version 0.0.2. Available at: <https://github.com/datalimited/datalimited>
- Anderson, S.C, A.B. Cooper, O.P. Jensen, C. Minto, J.T. Thorson, J.C. Walsh, J. Afflerbach, M. Dickey-Collas, K.M. Kleisner, C. Longo, G.C. Osio, D. Ovando, I. Mosqueira, A.A. Rosenberg, and E.R. Selig. 2017. Improving estimates of population status and trend with superensemble models. *Fish & Fisheries* 18(4):732-741.
- Berkson, J., L. Barbieri, S. Cadrin, S. Cass-Calay, P. Crone, M. Dorn, C. Friess, D. Kobayashi, T.J. Miller, W.S. Patrick, S. Pautzke, S. Ralston, and M. Trianni. 2011. Calculating Acceptable Biological Catch for stocks that have reliable catch data only (Only Reliable Catch Stocks – ORCS). NOAA Technical Memorandum NMFS-SEFSC-616. NMFS-SEFSC, Miami, FL.
- Carruthers, T.R., Punt, A.E., Walters, C.J., MacCall, A., McAllister, M.K., Dick, E.J. & Cope, J. 2014. Evaluating methods for setting catch limits in data-limited fisheries. *Fisheries Research*. 153: 48–68.
- Cohen, J. 1968. Weighted kappa: Nominal scale agreement provision for scaled disagreement or partial credit. *Psychological Bulletin* 70(4):213-220.
- Costello, C., D. Ovando, R. Hilborn, S.D. Gaines, O. Deschenes, and S.E. Lester. 2012. Status and solutions for the world's unassessed fisheries. *Science* 338:517-520.
- Csirke, J. and Sharp, G. (1984) Proceedings of the Expert Consultation to Examine Changes in Abundance and Species Composition of Neritic Fish Resources. FAO Fisheries Report No. 291, San Jose, Costa Rica, 102 p.
- Dick, E.J., and A.D. MacCall. 2011. Depletion-based stock reduction analysis: A catch-based method for determining sustainable yields for data-poor fish stocks. *Fisheries Research* 110:331-341.
- Fleiss, J.L. 1973. *Statistical methods for rates and proportions*. John Wiley & Sons, New York, NY.
- Free, C.M., O.P. Jensen, J. Wiedenmann, and J.J. Deroba. 2017. The refined ORCS approach: a catch-based method for estimating stock status and catch limits for data-poor fish stocks. *Fisheries Research* 193:60-70.

Jensen, O.P., Free, C.M. (2017) Testing and comparison of data-limited assessment models for estimating global and regional stock status. UN Food & Agriculture Organization.

Froese, R., and D. Pauly. 2016. FishBase. Available at: www.fishbase.org

Froese, R., and K. Kesner-Reyes. 2002. Impact of Fishing on the Abundance of Marine Species. ICES CM 2002/L:12.

Froese, R., N. Demirel, G. Coro, K.M. Kleisner, and H. Winker. 2017. Estimating fisheries reference points from catch and resilience. *Fish & Fisheries* 18(3):506-526.

Kleisner, K., D. Zeller, R. Froese, and D. Pauly. 2013. Using global catch data for inferences on the world's marine fisheries. *Fish & Fisheries* 14(3):293-311.

Kleisner, K., and D. Pauly. 2011. Stock-status plots of fisheries for Regional Seas, in: Christensen, V., S. Lai, M.L.D. Palomares, D. Zeller, D. Pauly (Eds.), *The State of Biodiversity and Fisheries in Regional Seas*. Fisheries Centre Research Reports 19(3). University of British Columbia, Vancouver, Canada, pp. 37-40.

Landis, J.R., and G.G. Koch. 1977. The measurement of observer agreement for categorical data. *Biometrics* 33(1):159-174.

MacCall, A.D. 2009. Depletion-corrected average catch: A simple formula for estimating sustainable yields in data-poor situations. *ICES Journal of Marine Science* 66:2267-2271.

Martell, S., and R. Froese. 2013. A simple method for estimating MSY from catch and resilience. *Fish & Fisheries* 14:504-514.

Pauly, D., Hilborn, R., and T.A. Branch. 2013. Fisheries: Does catch reflect abundance? *Nature* 494(7437):303-6. DOI: 10.1038/494303a

Pinsky, M.L., O.P. Jensen, D. Ricard, S.R. Palumbi. 2011. Unexpected patterns of fisheries collapse in the world's oceans. *Proceedings of the National Academy of Sciences* 108(20):8317-8322.

Punt, A.E., Butterworth, D.S., de Moor, C.L., De Oliveira, J.A.A. & Haddon, M. 2014. Management strategy evaluation: best practices. *Fish and Fisheries*. 17: 303–334.

R Core Team. 2016. R: A language and environment for statistical computing. R Foundation for Statistical Computing: Vienna, Austria. Available at: <http://www.R-project.org/>

Ricard, D., C. Minto, O.P. Jensen, and J.K. Baum. 2012. Examining the knowledge base and status of commercially exploited marine species with the RAM Legacy Stock Assessment Database. *Fish & Fisheries* 13(4):380-398.

- Rosenberg, A.A., M.J. Fogarty, A.B. Cooper, M. Dickey-Collas, E.A. Fulton, N.L. Gutiérrez, K.J.W. Hyde, K.M. Kleisner, C. Longo, C.V. Minto-Vera, C. Minto, I. Mosqueira, G.C. Osio, D. Ovando, E.R. Selig, J.T. Thorson, and Y. Ye. 2014. Developing new approaches to global stock status assessment and fishery production potential of the seas. FAO Fisheries and Aquaculture Circular, Rome, Italy.
- Rosenberg, A.A., Kleisner, K.M., Afflerbach, J., Anderson, S.C., Dickey-Collas, M., Cooper, A.B., Fogarty, M.J., Fulton, E.A., Gutiérrez, N.L., Hyde, K.J.W., Jardim, E., Jensen, O.P., Kristiansen, T., Longo, C., Minto-Vera, C.V., Minto, C., Mosqueira, I., Chato Osio, G., Ovando, D., Selig, E.R., Thorson, J.T., Walsh, J.C., Ye, Y. (2017) Applying a new ensemble approach to estimating stock status of marine fisheries around the world. Conservation Letters: doi: 10.1111/conl.12363
- SAFMC. 2012. SSC ORCS Workshop Report. South Atlantic Fishery Management Council-Scientific & Statistical Committee (SAFMC-SSC), North Charleston, SC.
- SAFMC. 2013. SSC ORCS Workshop II Report. South Atlantic Fishery Management Council-Scientific & Statistical Committee (SAFMC-SSC), North Charleston, SC.
- Thorson, J.T., C. Minto, C.V. Minto-Vera, K.M. Kleisner, and C. Longo. 2013. A new role for effort dynamics in the theory of harvested populations and data-poor stock assessment. Canadian Journal of Fisheries & Aquatic Sciences 70(12):1829-1844.
- Vasconcellos, M., and K. Cochrane. 2005. Overview of world status of data-limited fisheries: inferences from landings statistics. In: Kruse, G.H., V.F. Gallucci, D.E. Hay, R.I. Perry, R.M. Peterman, T.C. Shirley, P.D. Spencer, B. Wilson, and R. Woodby (Eds.): Fisheries Assessment and Management in Data-Limited Situations. Alaska Sea Grant College Program, University of Alaska Fairbanks, Fairbanks, AK, pp. 1–20.
- Zhou, S., A.E. Punt, Y. Ye, N. Ellis, C.M. Dichmont, M. Haddon, D.C. Smith, and A.D.M. Smith. 2017. Estimating stock depletion level from patterns of catch history. Fish & Fisheries 18(4):742-751.
- Zhou, S., Z. Chen, C.M. Dichmont, N. Ellis, M. Haddon, A.E. Punt, A.D.M. Smith, D.C. Smith, and Y. Ye. 2016. An optimised catch-only assessment method for data poor fisheries. In: Catch-based methods for data-poor fisheries. Report to FAO. CSIRO, Brisbane, Australia.

Tables & Figures

Tables

Table 1. Catch-only stock assessment models

Table 2. Superensemble model parameters and goodness of fit statistics

Figures

Figure 1. Overall performance (continuous/categorical) – simulated and real stocks

Figure 2. Superensemble model performance (observed vs. predicted)

Figure 3. Performance by ED scenario (continuous/categorical) – simulated stocks

Figure 4. Performance by ED scenario (continuous/categorical) – real stocks

Figure 5. Performance by LH-ED scenario (continuous) – simulated stocks

Figure 6. Performance by LH-ED scenario (categorical) – simulated stocks

Figure 7. Performance by LH-ED scenario (continuous) – real stocks

Figure 8. Performance by LH-ED scenario (categorical) – real stocks

Figure 9. ED misspecification consequences (continuous/categorical) – simulated stocks

Figure 10. LH-ED misspecification consequences (continuous) – simulated stocks

Figure 11. LH-ED misspecification consequences (categorical) – simulated stocks

Figure 12. cMSY-17 and 2-Stage saturation priors

Figure 13. Performance of cMSY-17 and 2-Stage saturation priors – simulated/real stocks

Figure 14. Two-stage catch only model performance (observed vs. predicted)

Supporting Tables & Figures

Supp. Table 1. Simulated stock factorial design

Supp. Table 2. Simulated stock life history categories

Supp. Table 3. Simulated stock effort dynamics categories

Supp. Table 4. RAMLDB stock life history classification key

Supp. Figure 1. Overall BRT model tuning

Supp. Figure 2. ED BRT model tuning

Appendices

Appendix A. LH-ED BRT model tuning

Appendix B. RAMLDB effort dynamics classification

Tables & Figures

Table 1. Catch-only stock assessment methods.

	Method	References	Data input/output	Brief description
1	rORCS Refined ORCS approach	Berkson et al. 2011 Free et al. 2017	In: Catch, 12 questions Out: Exploitation status, OFL	Uses a boosted classification tree model trained on the RAMLDB to predict status from 12 stock- and fishery-related predictors
2	cMSY-2013 Catch-MSY	Martell & Froese 2013 Rosenberg et al. 2014	In: Catch, resilience Out: B/B_{MSY} , MSY, B , B_{MSY}	Uses a stock reduction analysis with priors for r , k , and initial/final year depletion derived from resilience to estimate status
3	cMSY-2017*† Updated catch-MSY	Froese et al. 2017	In: Catch, resilience Out: B/B_{MSY} , all MSY ref points	Updates the cMSY-2013 stock reduction analysis with a new algorithm for identifying probable r - k pairs to estimate status
4	COM-SIR* Catch-only-model with sampling importance resampling	Vasconcellos & Cochrane 2005 Rosenberg et al. 2014	In: Catch, resilience Out: B/B_{MSY}	Uses a coupled harvest-dynamics model fit using a sampling importance resampling algorithm to estimate status
5	SSCOM* State-space catch-only model	Thorson et al. 2013	In: Catch, resilience Out: B/B_{MSY}	Uses a coupled harvest-dynamics model fit using a Bayesian hierarchical state-space framework to estimate status
6	SSP-2002 Stock status plots	Froese & Kesner-Reyes 2002	In: Catch Out: Development status	Uses simple rules that compare the final year's catch to the maximum catch to estimate status
7	SSP-2013 Updated stock status plots	Kleisner et al. 2013 Kleisner & Pauly 2011	In: Catch Out: Development status	Updates the rules of SSP-2002 to utilize the minimum catch occurring after the maximum catch to estimate status
8	mPRM* Modified panel regression model	Costello et al. 2012 Anderson et al. 2017	In: Catch, taxonomic group Out: B/B_{MSY}	Uses a panel regression model trained on the RAMLDB to predict status from characteristics of the catch time series and taxonomic group
9	Zhou-BRT*† Catch-only boosted regression trees	Zhou et al. 2017	In: Catch Out: Saturation	Uses a boosted regression tree model trained on the RAMLDB to predict status from 38 catch history statistics
10	Zhou-OCOM Optimized catch-only model	Zhou et al. 2016	In: Catch, natural mortality (M) Out: Saturation, MSY	Uses a stock reduction analysis with priors for r and final year depletion derived from M and saturation from Zhou-BRT to estimate status

* used to develop the superensemble models

† coupled to create the two-stage catch-only model

Table 2. Optimal BRT model parameters and cross-validation goodness of fit statistics for the superensemble models.

Model	# of trees	Interaction depth	Learning rate	RMSE	r²
Overall model (n=1)	7500	10	0.005	0.50	0.37
ED models (n=4)					
Constant F	5000	12	0.001	0.32	0.21
Biomass-coupled F	4000	12	0.005	0.27	0.76
Increasing F	8500	12	0.001	0.26	0.60
Rollercoaster F	8000	12	0.01	0.21	0.62
LH-ED models (n=12)					
<i>Small pelagic</i>					
Constant F	3500	6	0.001	0.27	0.17
Biomass-coupled F	4000	12	0.001	0.20	0.31
Increasing F	6000	10	0.001	0.27	0.53
Rollercoaster F	10000	12	0.005	0.21	0.67
<i>Large pelagic</i>					
Constant F	100	12	0.005	0.33	0.03
Biomass-coupled F	4500	10	0.001	0.19	0.91
Increasing F	3500	12	0.001	0.22	0.46
Rollercoaster F	10000	10	0.01	0.16	0.52
<i>Demersal</i>					
Constant F	3000	12	0.001	0.31	0.32
Biomass-coupled F	4000	12	0.005	0.34	0.73
Increasing F	6000	12	0.001	0.26	0.65
Rollercoaster F	7500	6	0.005	0.21	0.71

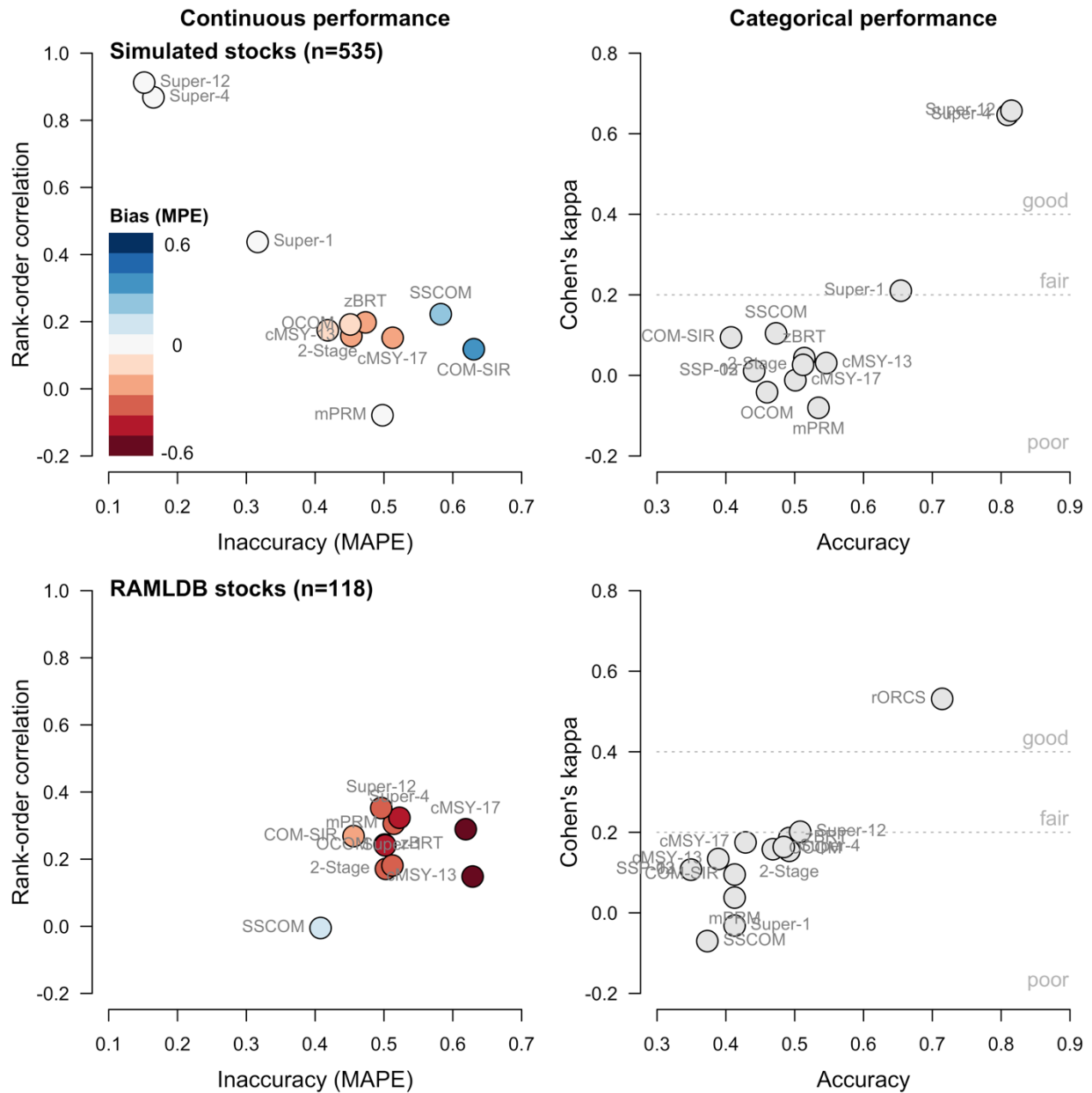


Figure 1. The continuous and categorical performance of COMs evaluated on the simulated stocks withheld from the BRT model training (n=535) and (2) RAMLDB stocks (n=118). In the continuous performance plots, the best performing methods are indicated by high rank-order correlation and low inaccuracy (top-left corner). In the categorical performance plots, the best performing methods are indicated by high Cohen's kappa and high accuracy (top-right corner).

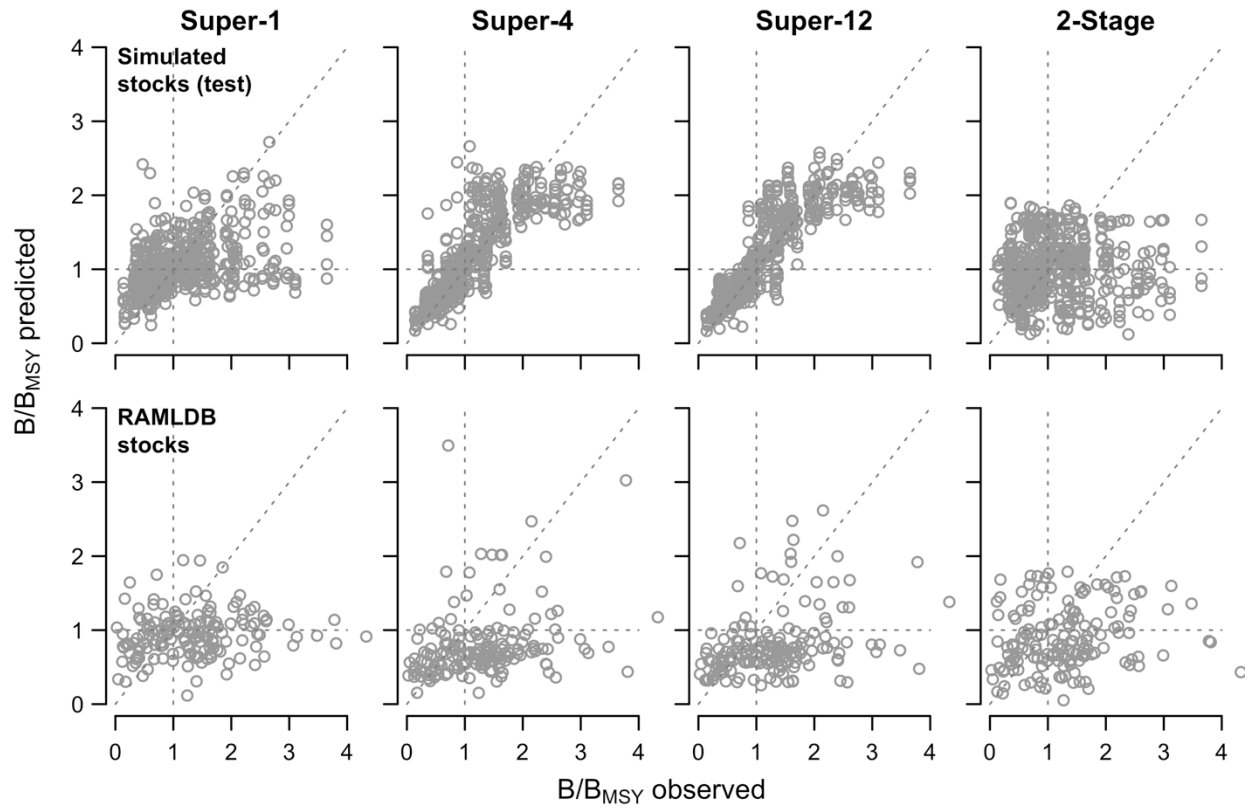


Figure 2. Observed stock status versus stock status predicted by the new superensemble and two-stage models for the test simulated stocks (n=576) and RAMLDB stocks (n=161).

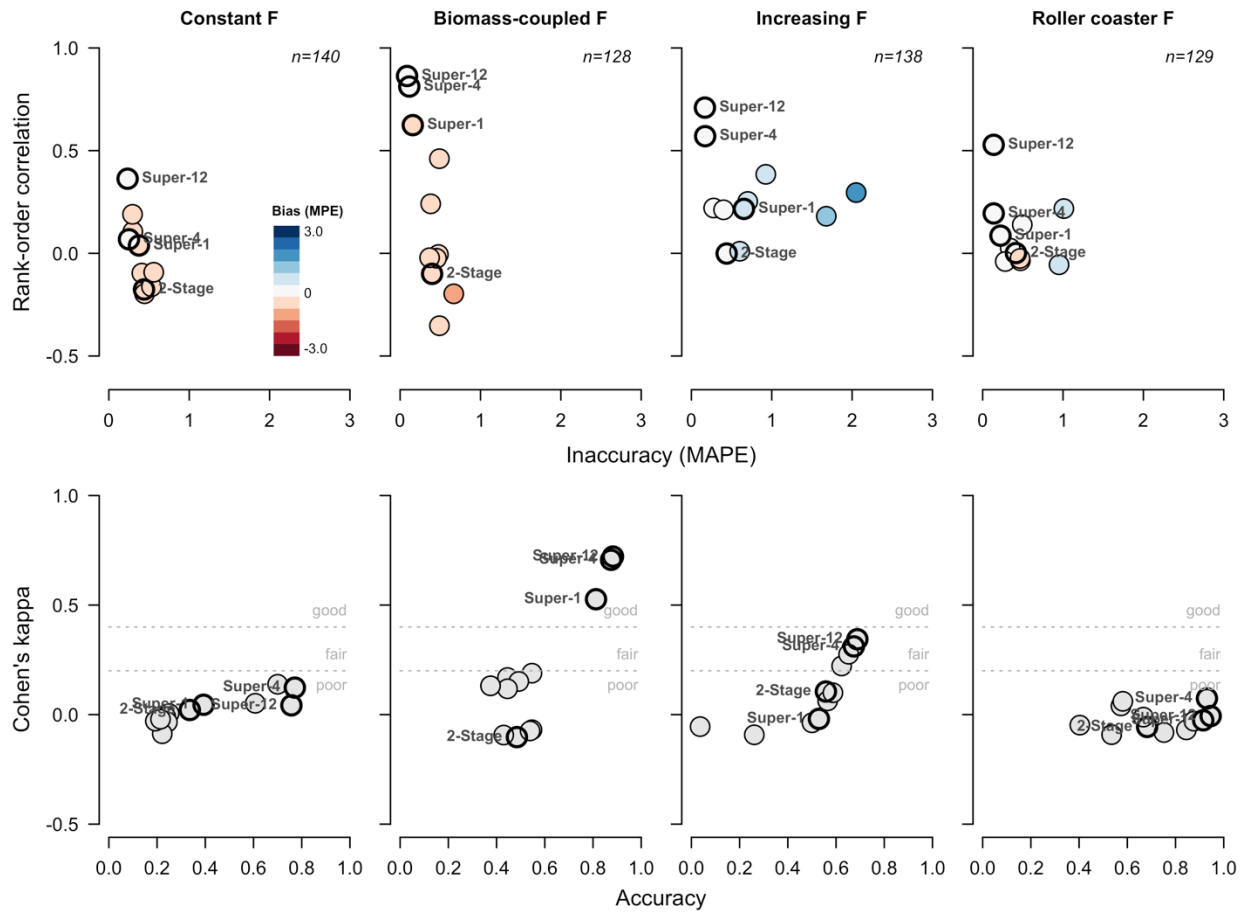


Figure 3. The continuous and categorical performance of COMs evaluated on the test simulated stocks by effort dynamics (144 stock max per scenario). In the continuous performance plots, the best performing methods are indicated by high rank-order correlation and low inaccuracy (top-left corner). In the categorical performance plots, the best performing methods are indicated by high Cohen's kappa and high accuracy (top-right corner). Samples sizes are shown in the top-right corner of each plot and are sometimes less than the 144 stock maximum due to failed convergence by one or more of the evaluated COMs.

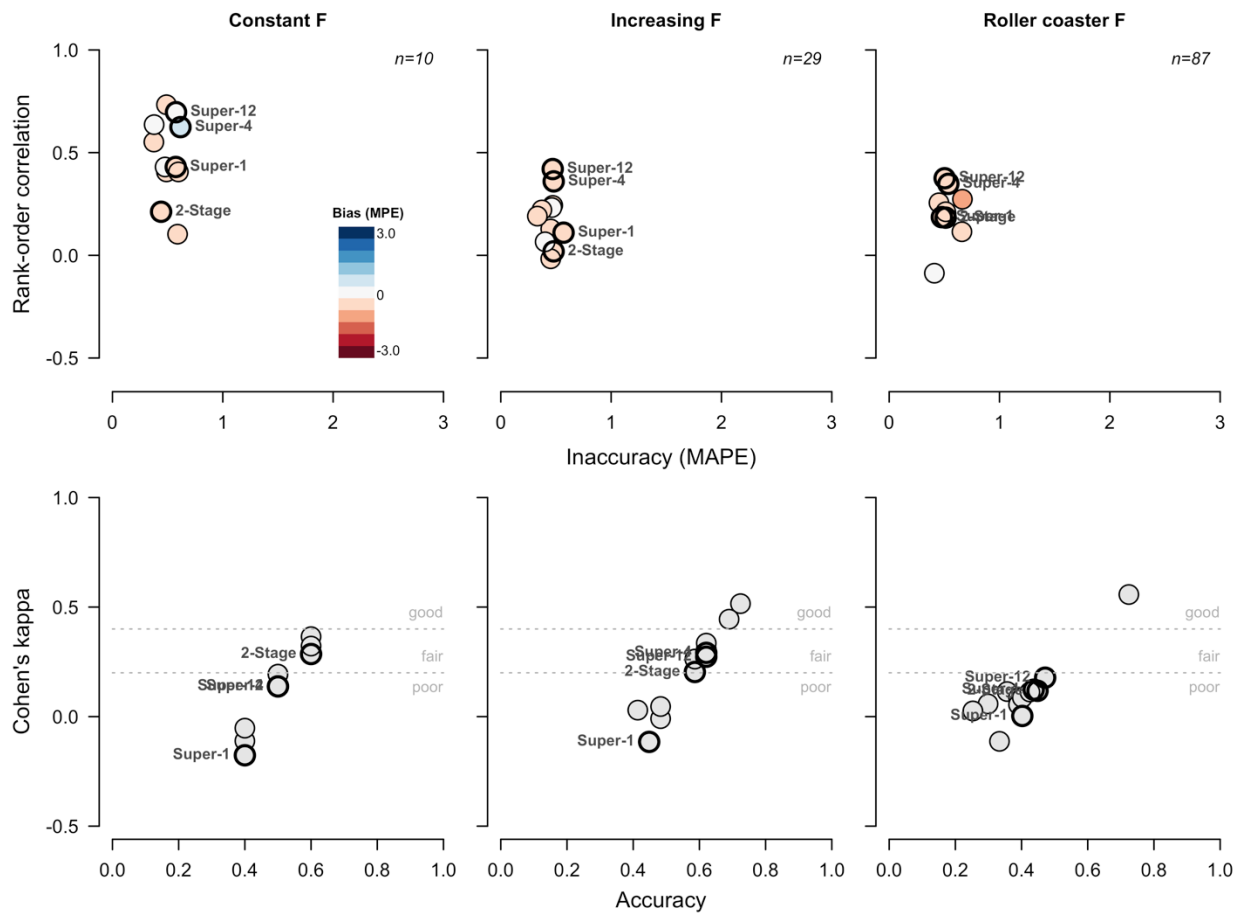


Figure 4. The continuous and categorical performance of COMs evaluated on the RAMLDB stocks by effort dynamics (note: sample sizes are not uniform). In the continuous performance plots, the best performing methods are indicated by high rank-order correlation and low inaccuracy (top-left corner). In the categorical performance plots, the best performing methods are indicated by high Cohen's kappa and high accuracy (top-right corner). Samples sizes are shown in the top-right corner of each plot and are sometimes less than the maximum possible due to failed convergence by one or more of the evaluated COMs.

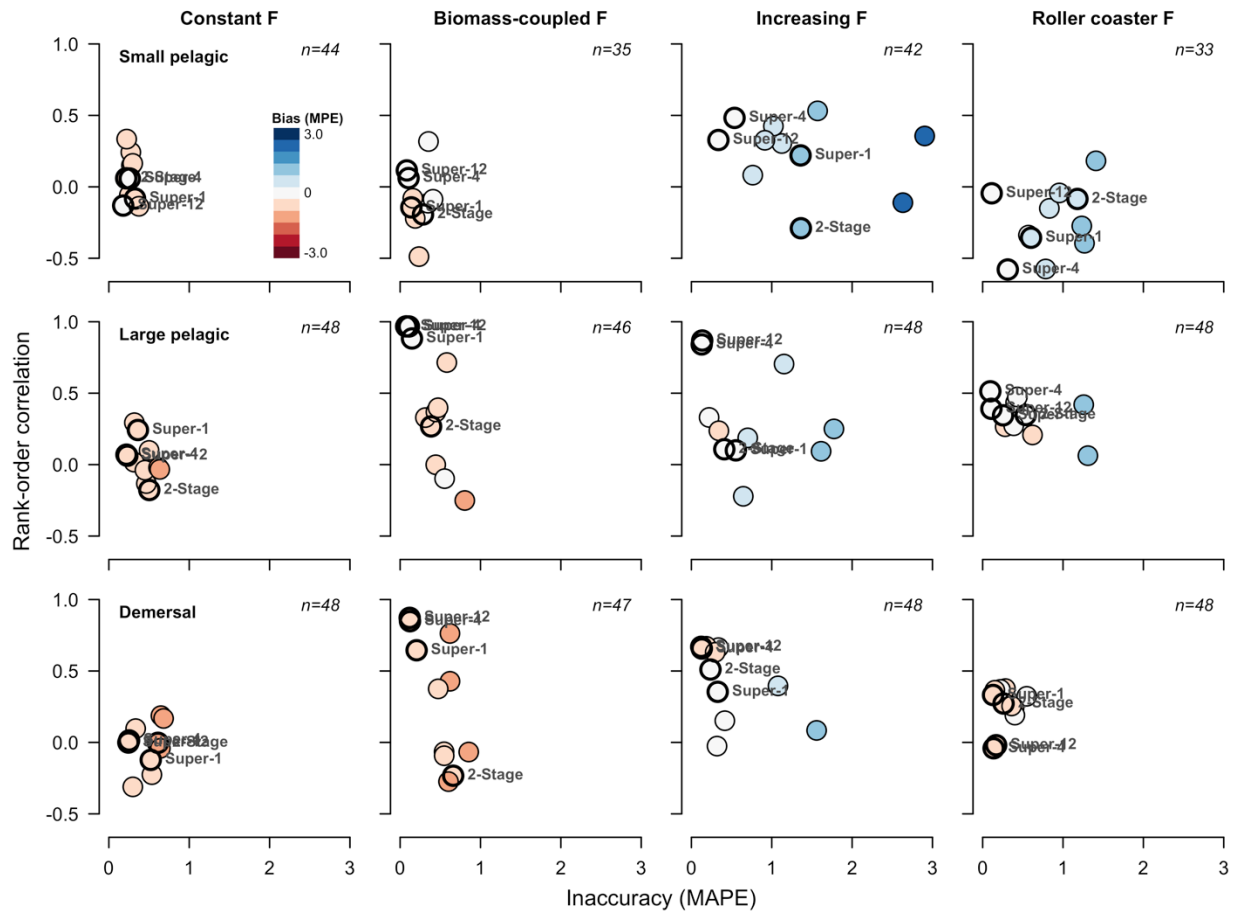


Figure 5. The continuous performance of COMs evaluated on the test simulated stocks by LH-ED couple (48 stock max per scenario). The best performing methods are indicated by high rank-order correlation and low inaccuracy (top-left corner). To maximize clarity, only the superensemble and two-stage catch-only models are labelled. Samples sizes are shown in the top-right corner of each plot and are sometimes less than the 48 stock maximum due to failed convergence by one or more of the evaluated COMs.

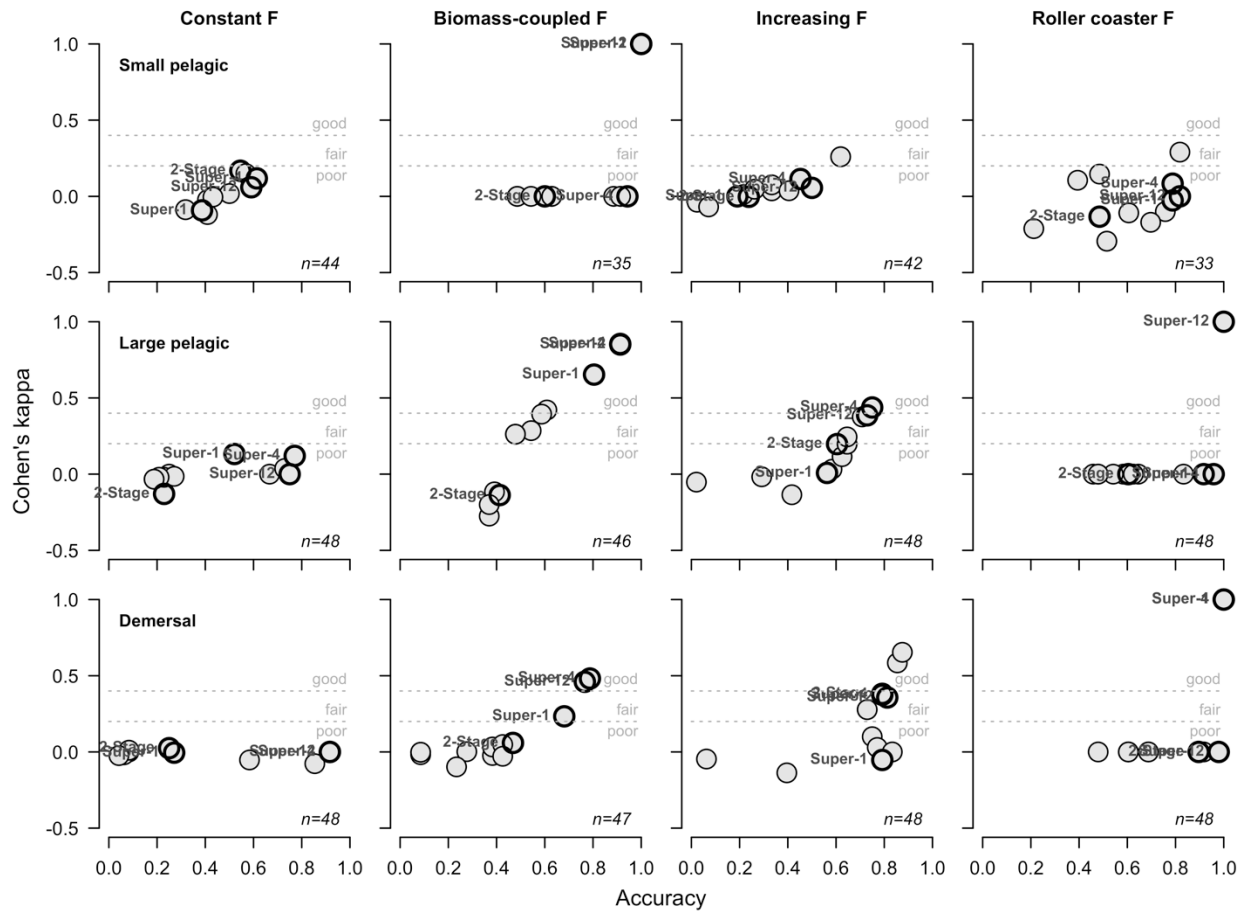


Figure 6. The categorical performance of COMs evaluated on the test simulated stocks by LH-ED couple (48 stock max per scenario). The best performing methods are indicated by high Cohen's kappa and high accuracy (top-right corner). To maximize clarity, only the superensemble and two-stage catch-only models are labelled. Samples sizes are shown in the bottom-right corner of each plot and are sometimes less than the 48 stock maximum due to failed convergence by one or more of the evaluated COMs. Note: Cohen's kappa cannot be estimated for classifiers that are 100% accurate (i.e., Cohen's kappa can never equal 1.0); however, classifiers achieving 100% accuracy were awarded a Cohen's kappa of 1.0 to allow visualization.

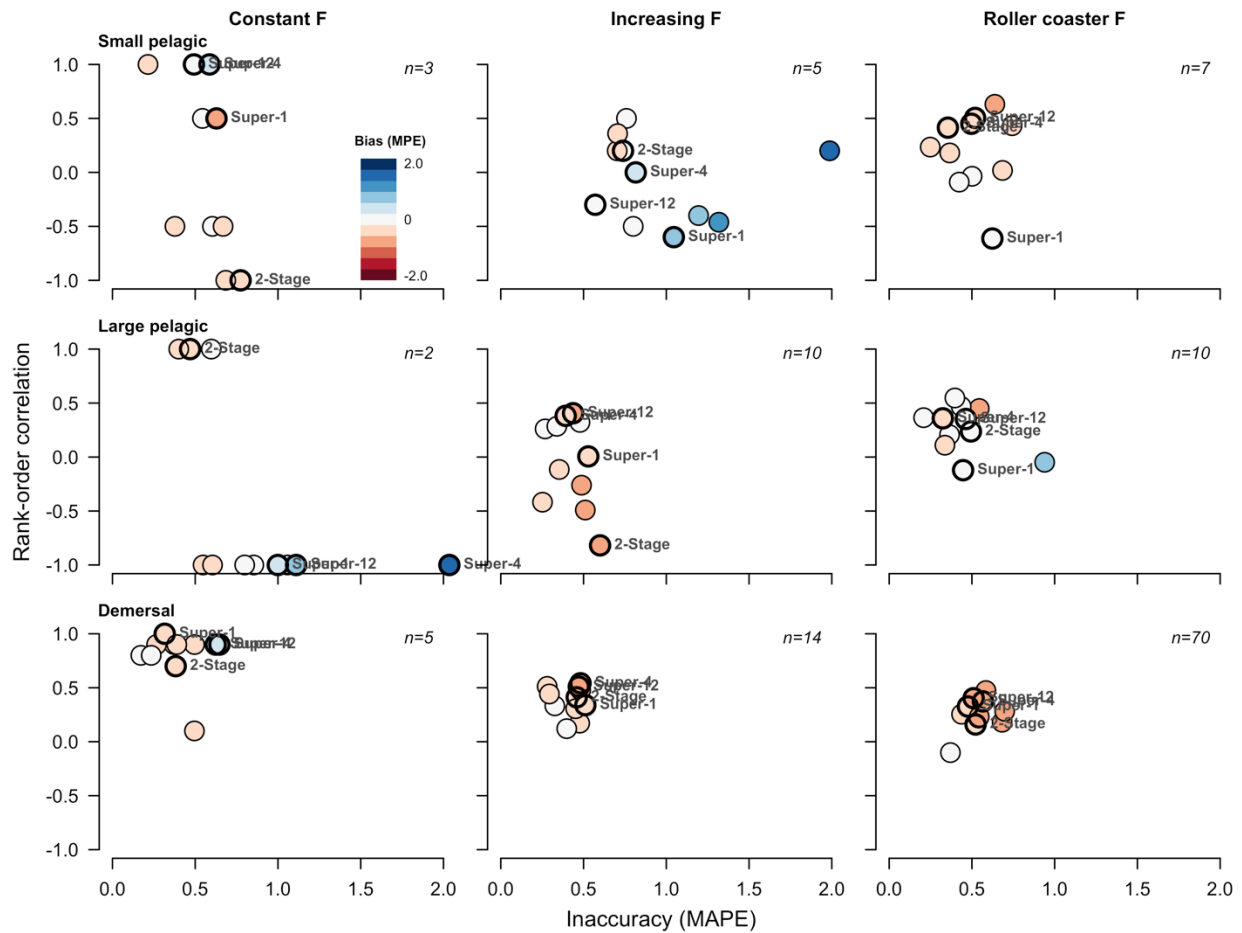


Figure 7. The continuous performance of COMs evaluated on the RAMLDB stocks by LH-ED couple (note: sample sizes are not uniform). The best performing methods are indicated by high rank-order correlation and low inaccuracy (top-left corner). To maximize clarity, only the superensemble and two-stage catch-only models are labelled. Sample sizes are shown in the top-right corner of each plot and are sometimes less than the maximum possible due to failed convergence by one or more of the evaluated COMs.

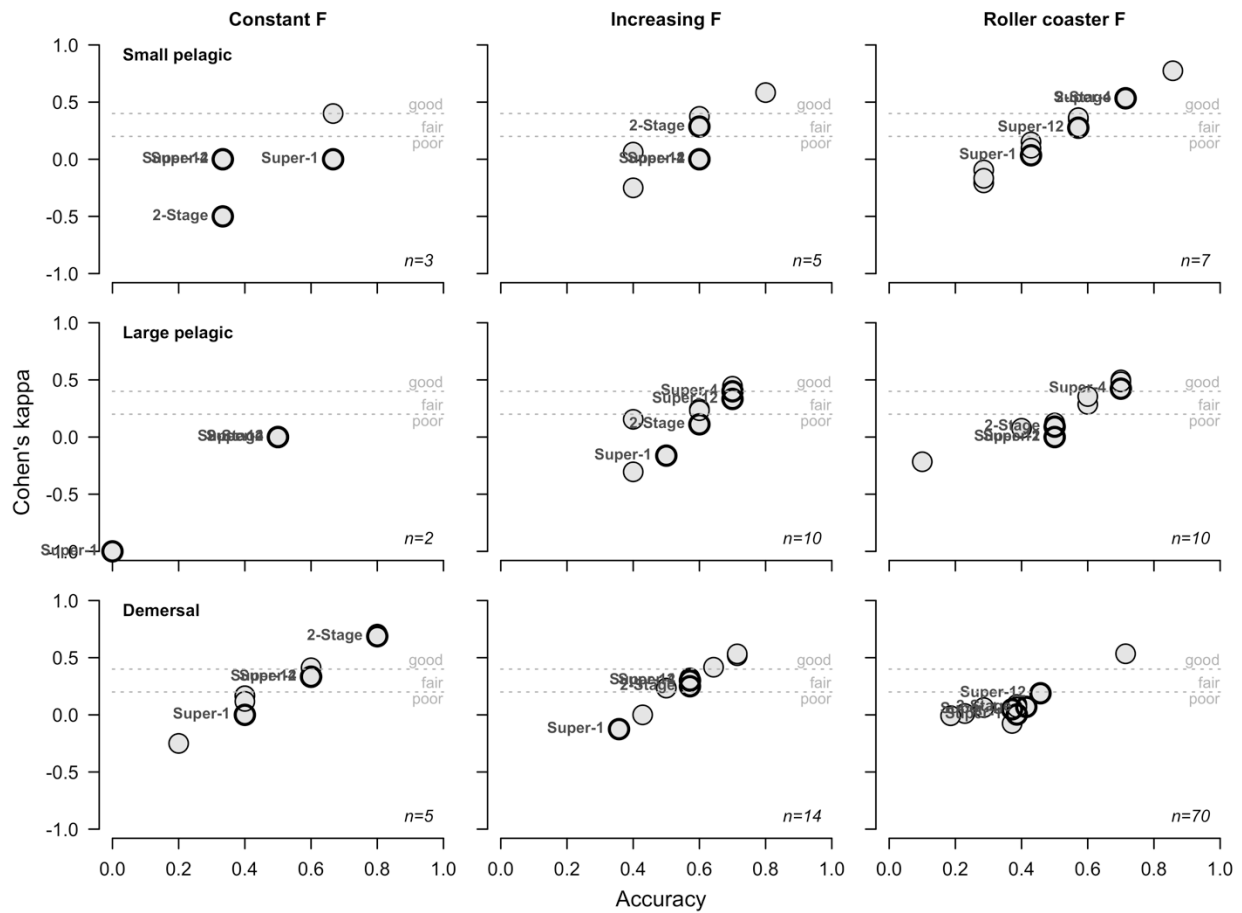


Figure 8. The categorical performance of COMs evaluated on the RAMLDB stocks by LH-ED couple (note: sample sizes are not uniform). The best performing methods are indicated by high Cohen's kappa and high accuracy (top-right corner). To maximize clarity, only the superensemble and two-stage catch-only models are labelled. Samples sizes are shown in the bottom-right corner of each plot and are sometimes less than the maximum possible due to failed convergence by one or more of the evaluated COMs.

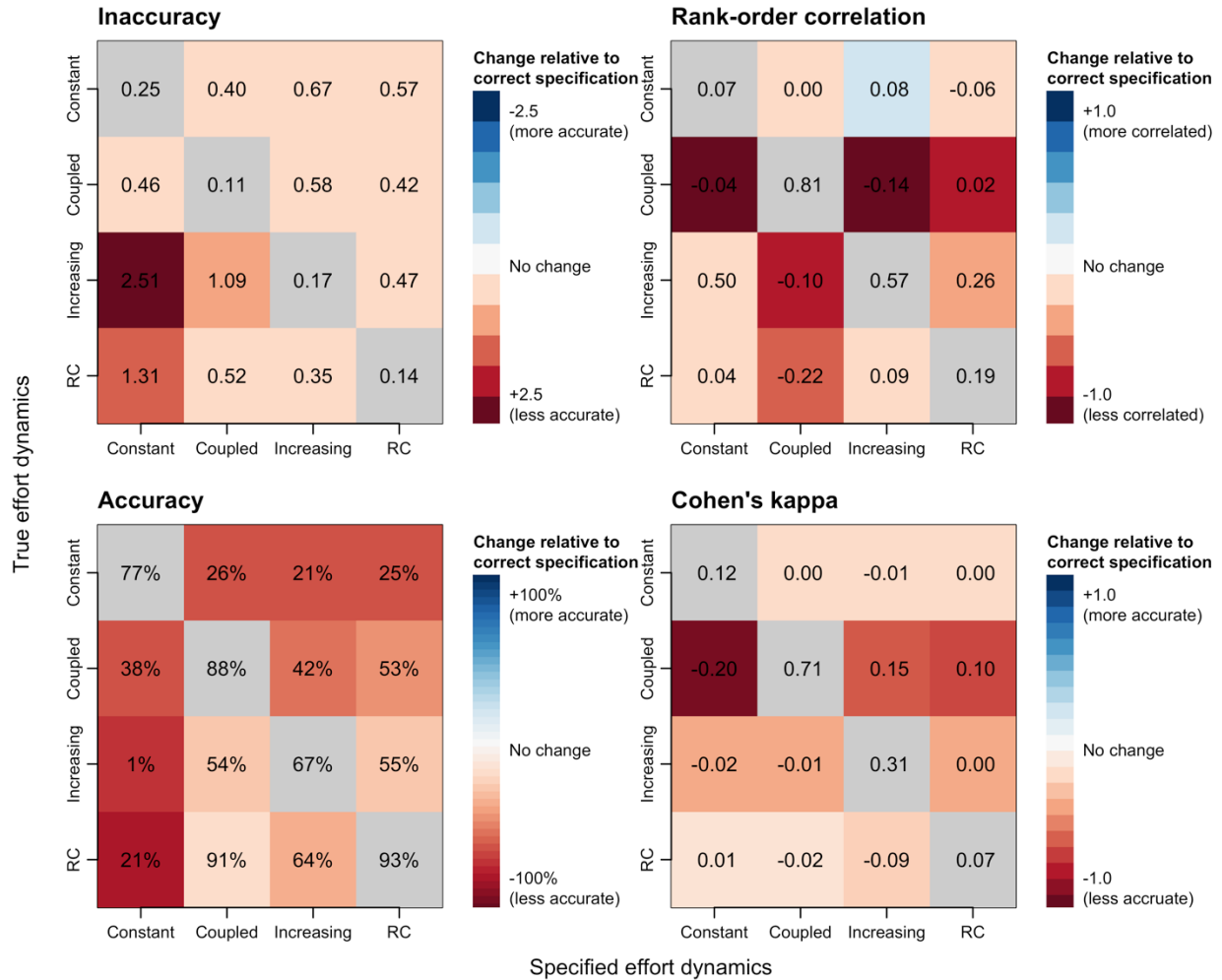


Figure 9. The sensitivity of continuous (inaccuracy and rank-order correlation) and categorical (accuracy and Cohen’s kappa) predictive performance of the ED superensemble models to the misspecification of the effort dynamics experienced by the test simulated stocks. The grey cells along the diagonal from top-left to bottom-right show the performance metric when effort dynamics are correctly specified. The within-row, off-diagonal cells show the performance metric when effort dynamics are incorrectly specified and the shading color indicates the change in performance relative to the correctly specified model (red=performance weakens, blue=performance improves, white=performance does not change).

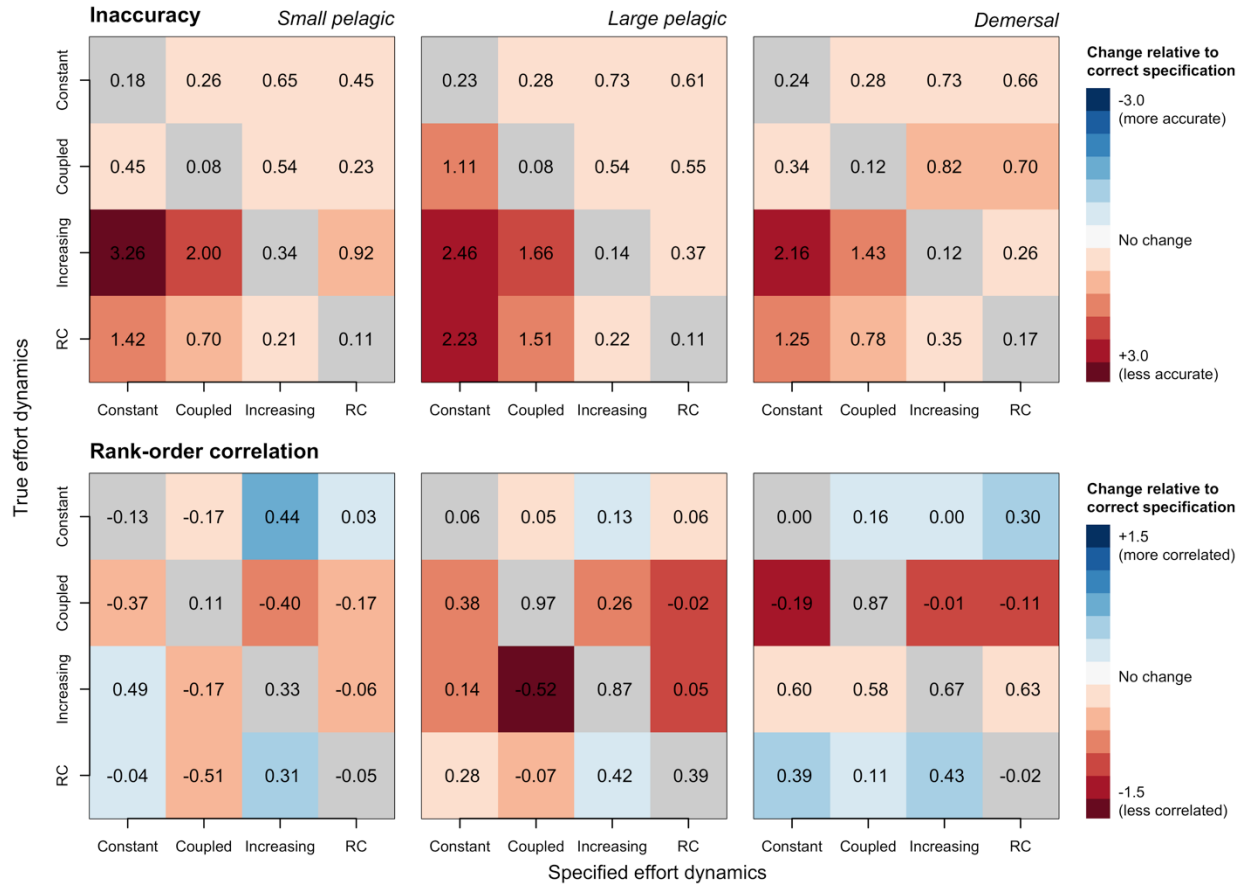


Figure 10. The sensitivity of continuous predictive performance (inaccuracy and rank-order correlation) of the LH-ED superensemble models to the misspecification of the effort dynamics experienced by the test simulated stocks. The grey cells along the diagonal from top-left to bottom-right show the performance metric when effort dynamics are correctly specified. The within-row, off-diagonal cells show the performance metric when effort dynamics are incorrectly specified and the shading color indicates the change in performance relative to the correctly specified model (red=performance weakens, blue=performance improves, white=performance does not change). In these comparisons, the life history category was always correctly specified.

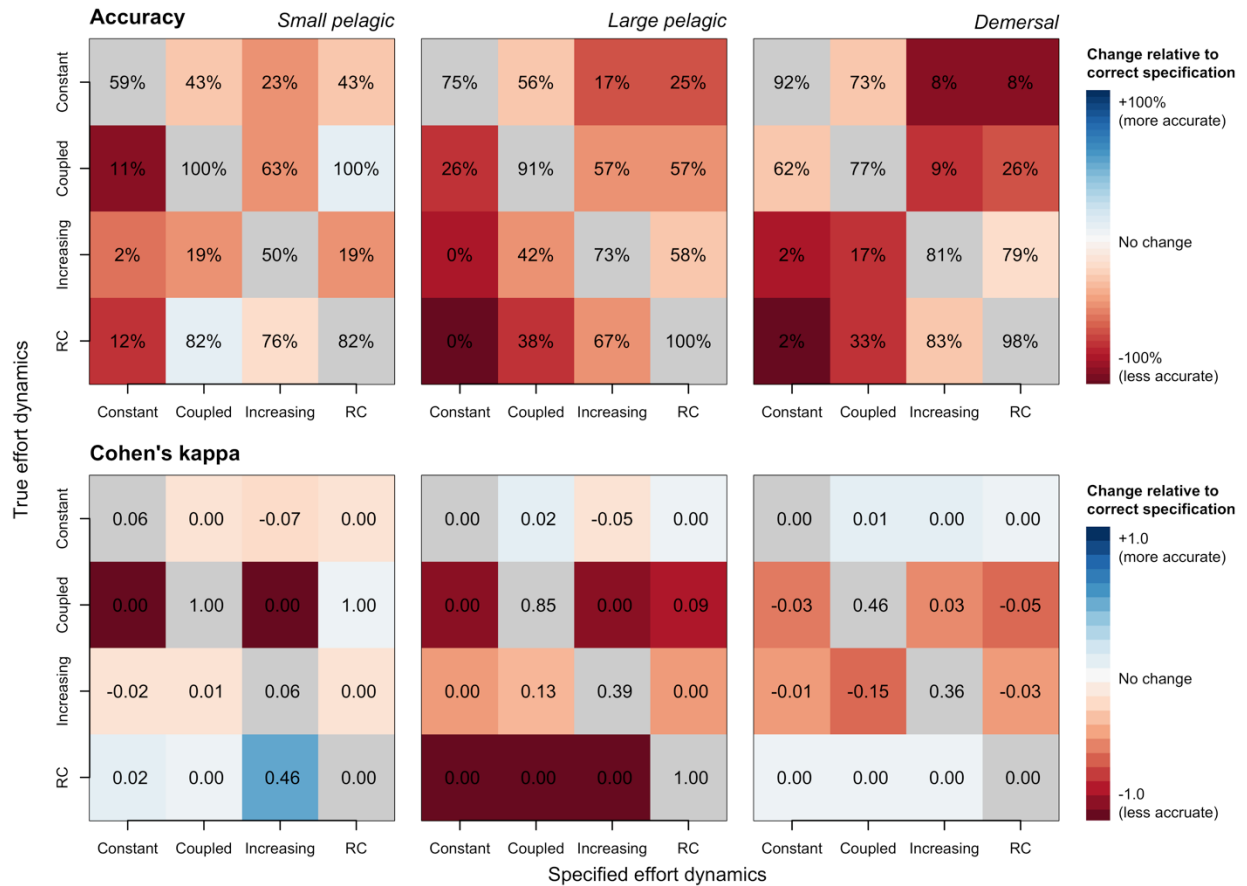


Figure 11. The sensitivity of categorical predictive performance (accuracy and Cohen's kappa) of the LH-ED superensemble models to the misspecification of the effort dynamics experienced by the test simulated stocks. The grey cells along the diagonal from top-left to bottom-right show the performance metric when effort dynamics are correctly specified. The within-row, off-diagonals cells show the performance metric when effort dynamics are incorrectly specified and the shading color indicates the change in performance relative to the correctly specified model (red=performance weakens, blue=performance improves, white=performance does not change). In these comparisons, the life history category was always correctly specified.

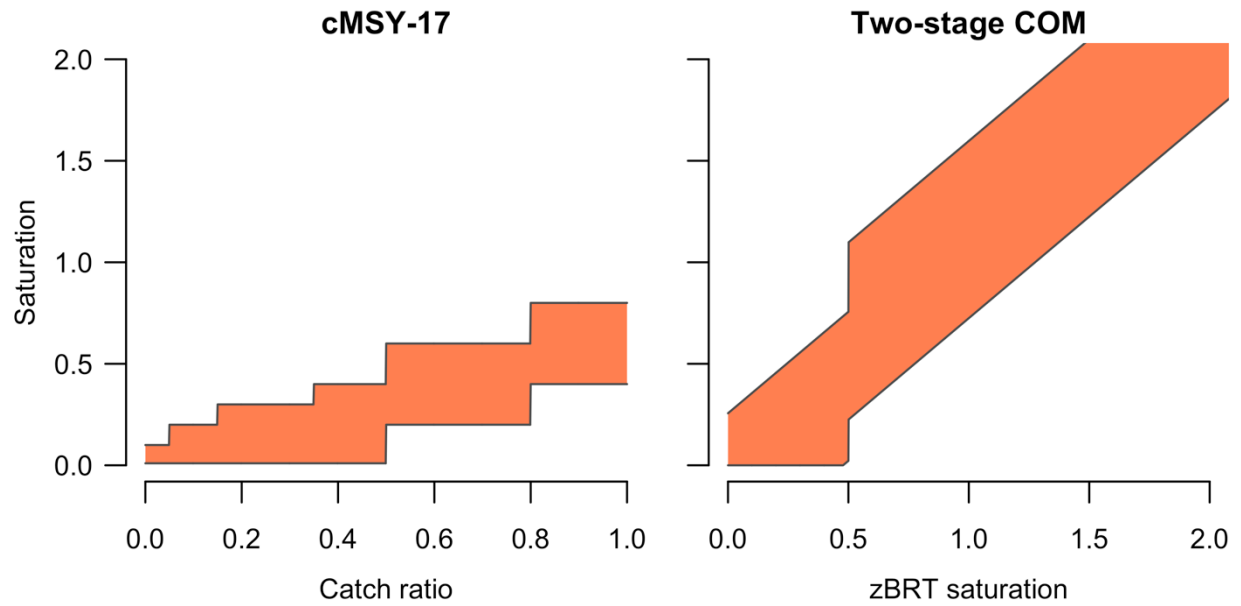


Figure 12. Visualization of the uniform final year saturation priors used by cMSY-17 and the new two-stage catch-only model. The cMSY-17 saturation priors are derived from catch ratio (last catch / maximum catch) in the final year while the two-stage model priors are derived from the zBRT saturation predictions in the final year.

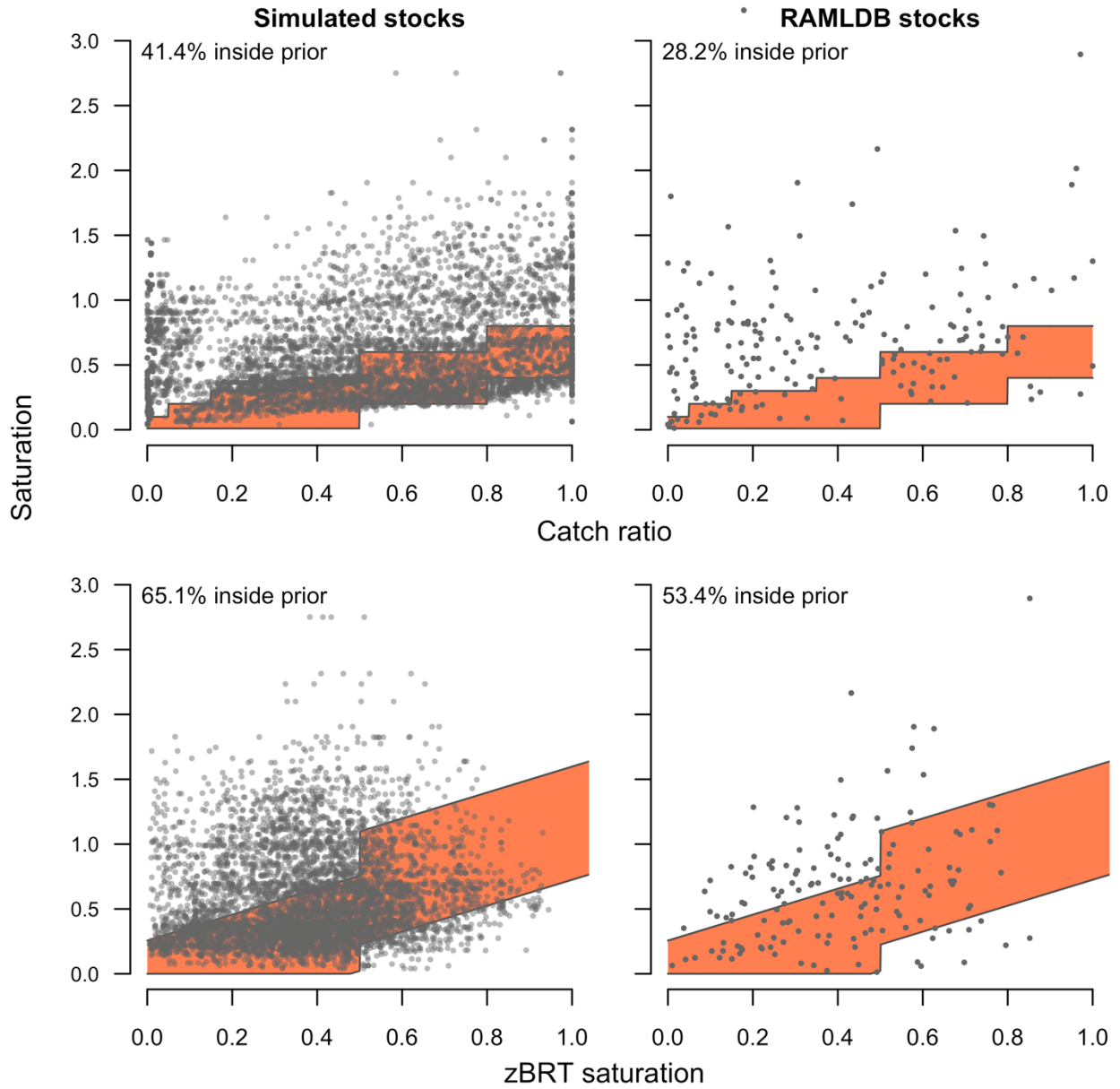


Figure 13. The suitability of the saturation priors used by cMSY-17 and the two-stage catch-only model for describing terminal year saturation for the simulated and RAMLDB stocks. The percentage of terminal year saturations contained by the prior is shown in the top-left corner.

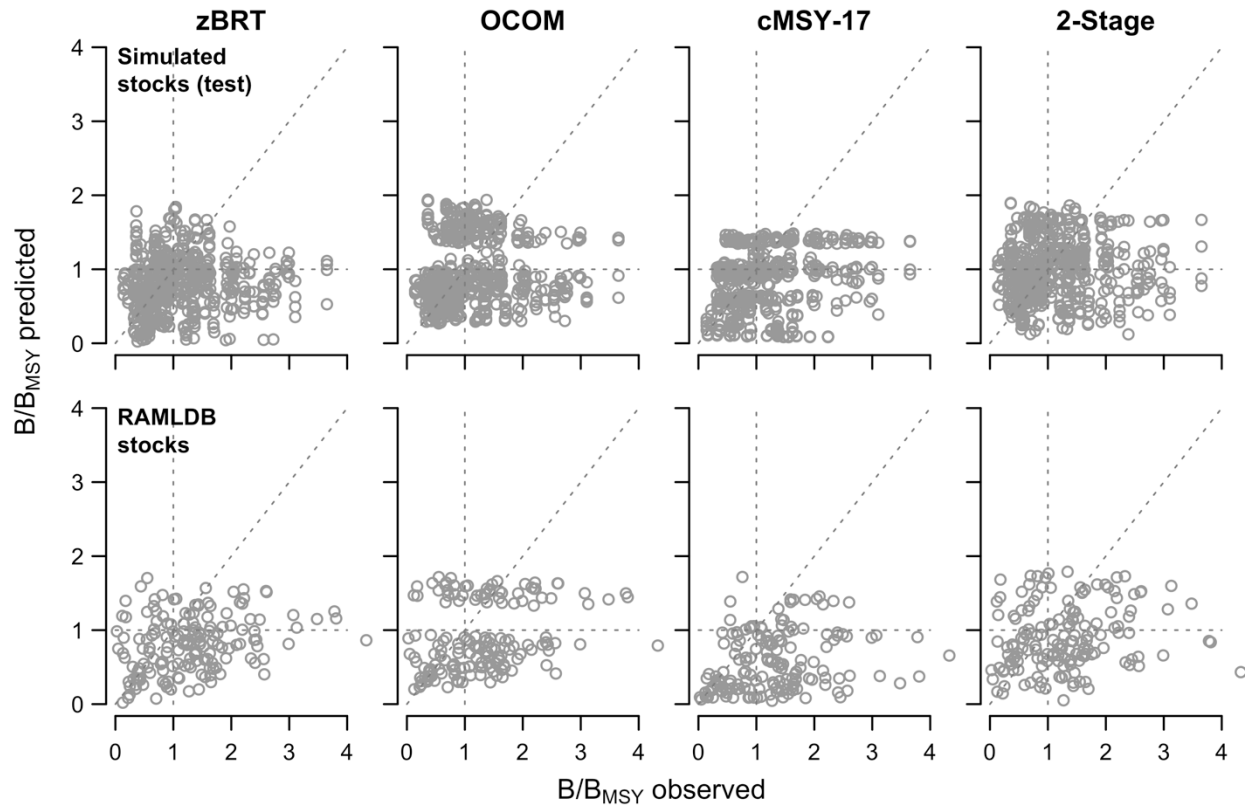


Figure 14. Observed stock status versus stock status predicted by four catch-only models for the test simulated stocks ($n=576$) and RAMLDB stocks ($n=161$). The two-stage model was created by coupling zBRT and cMSY-17 and OCOM is similar in structure.

Supporting Tables & Figures

Supp. Table 1. Factorial design of the Rosenberg et al. (2014) simulated stocks.

Factor	# of levels	Levels
Life history	3	Demersal, small pelagic, or large pelagic
Initial biomass depletion	3	100%, 70%, or 40% of carrying capacity
Exploitation dynamics	4	Constant, biomass-coupled, increasing, or roller coaster rates
Recruitment variability	2	Low or high variability
Recruitment autocorrelation	2	With or without autocorrelation
Catch measurement error	2	With or without catch measurement error
Time series length	2	20 or 60 years
Iterations	10	Iterations for each combination of the above parameters
Total # of stocks:	5760	

Supp. Table 2. Resilience and natural mortality (M) values for the life histories represented in the Rosenberg et al. (2014) simulated stocks.

Life history category	Generic name	Resilience	L_{inf} (cm)	T_{max} (yr)	M (yr⁻¹)*
Demersal	Gadoid	low	70	20	0.315
Small pelagic	Clupeoid	medium	30	8	0.729
Large pelagic	Scombrid	low	150	20	0.315

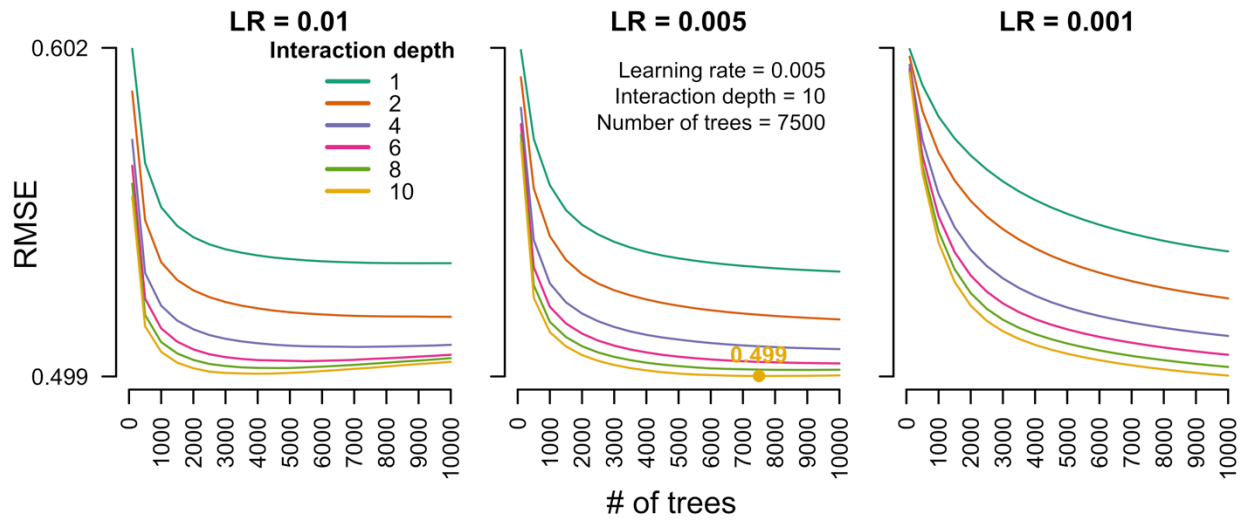
* Estimated using the t_{max}-based Hoenig (1983) method: $M = 4.899 * t_{max}^{-0.916}$

Supp. Table 3. Effort dynamics scenarios driving population dynamics in the Rosenberg et al. (2014) simulated stocks.

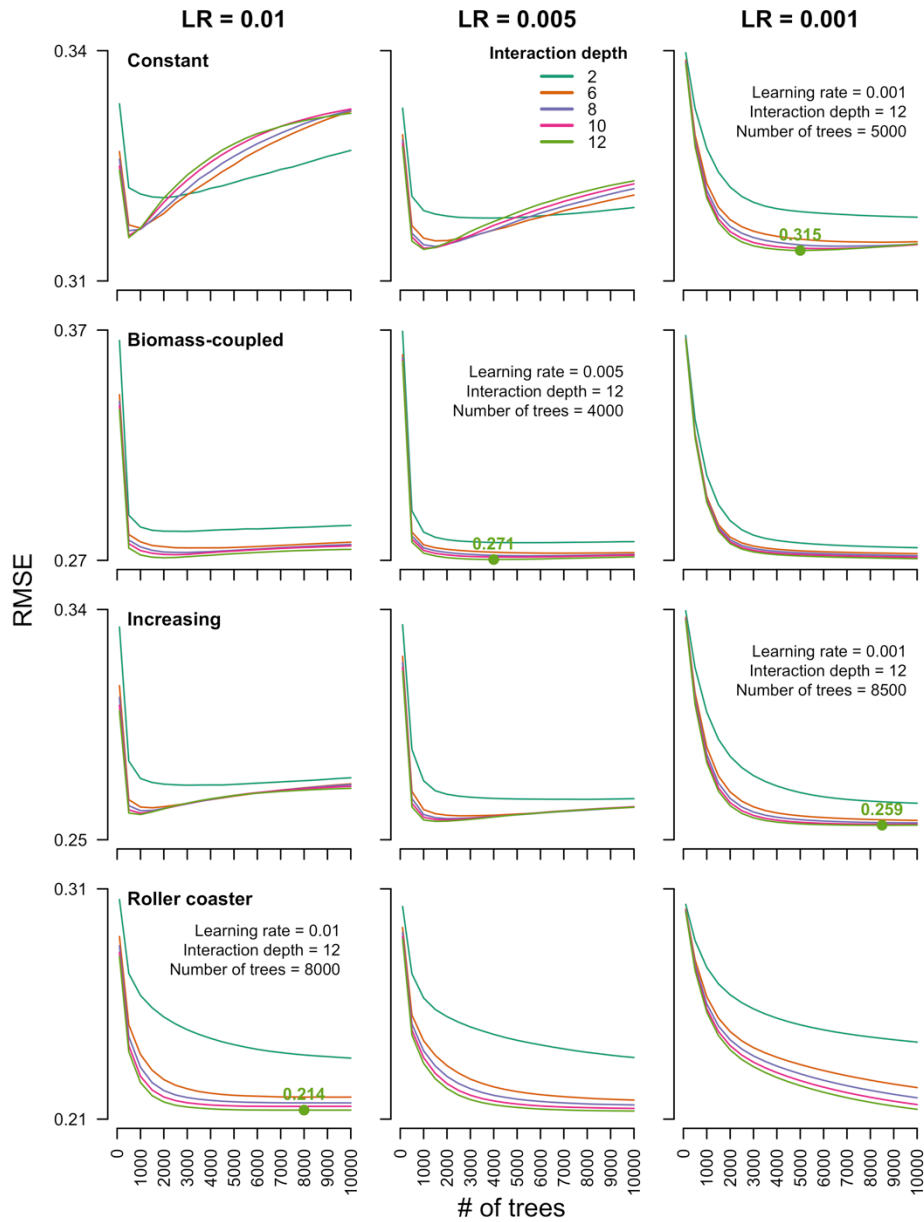
Effort dynamics scenario	Code	# of stocks		Description
		Simulated	RAMLDB	
Constant	ED0	1440	11	Harvest rate remains constant irrespective of biomass (e.g., bycatch species harvest)
Biomass-coupled	ED0.6	1440	---	Harvest rate has a dynamic relationship with biomass, see Rosenberg et al. (2014) for more details (e.g., open-access single-species harvest)
Increasing (one-way trip)	OW	1440	41	Harvest rate increases 5% per year to 80% of the harvest rate at which the stock crashes (e.g., a stock where harvest rate has continually increased)
Roller coaster (dome-shaped)	RC	1440	109	Harvest rate increases 25% per year to 80% of harvest rate at which the stock crashes, stays at this level for five years, then decreases to F_{MSY} levels by 30% per year (e.g., a stock where management began following extensive depletion)

Supp. Table 4. Mapping RAMLDB taxonomic families to the life histories represented in the simulated stocks.

Life history	# of stocks	Families
Demersal	116	Anoplopomatidae, Cheilodactylidae, Cottidae, Epigonidae, Gadidae, Hexagrammidae, Lutjanidae, Malacanthidae, Merlucciidae, Merlucciinae, Ophidiidae, Oreosomatidae, Paralichthyidae, Platycephalidae, Pleuronectidae, Rajidae, Scorpaenidae, Serranidae, Sparidae, Trachichthyidae, Uranoscopidae
Large pelagic	27	Istiophoridae, Pomatomidae, Sciaenidae, Scombridae, Xiphiidae
Small pelagic	18	Arripidae, Carangidae, Centrolophidae, Clupeidae, Engraulidae, Gempylidae, Sillaginidae, Stromateidae



Supp. Figure 1. Model tuning curves showing the average root mean square error (RMSE) for each combination of candidate BRT model parameters (learning rate, interaction depth, # of trees) for the overall superensemble model. The optimal combination of model parameters (marked and labeled) is the combination that minimizes the RMSE.



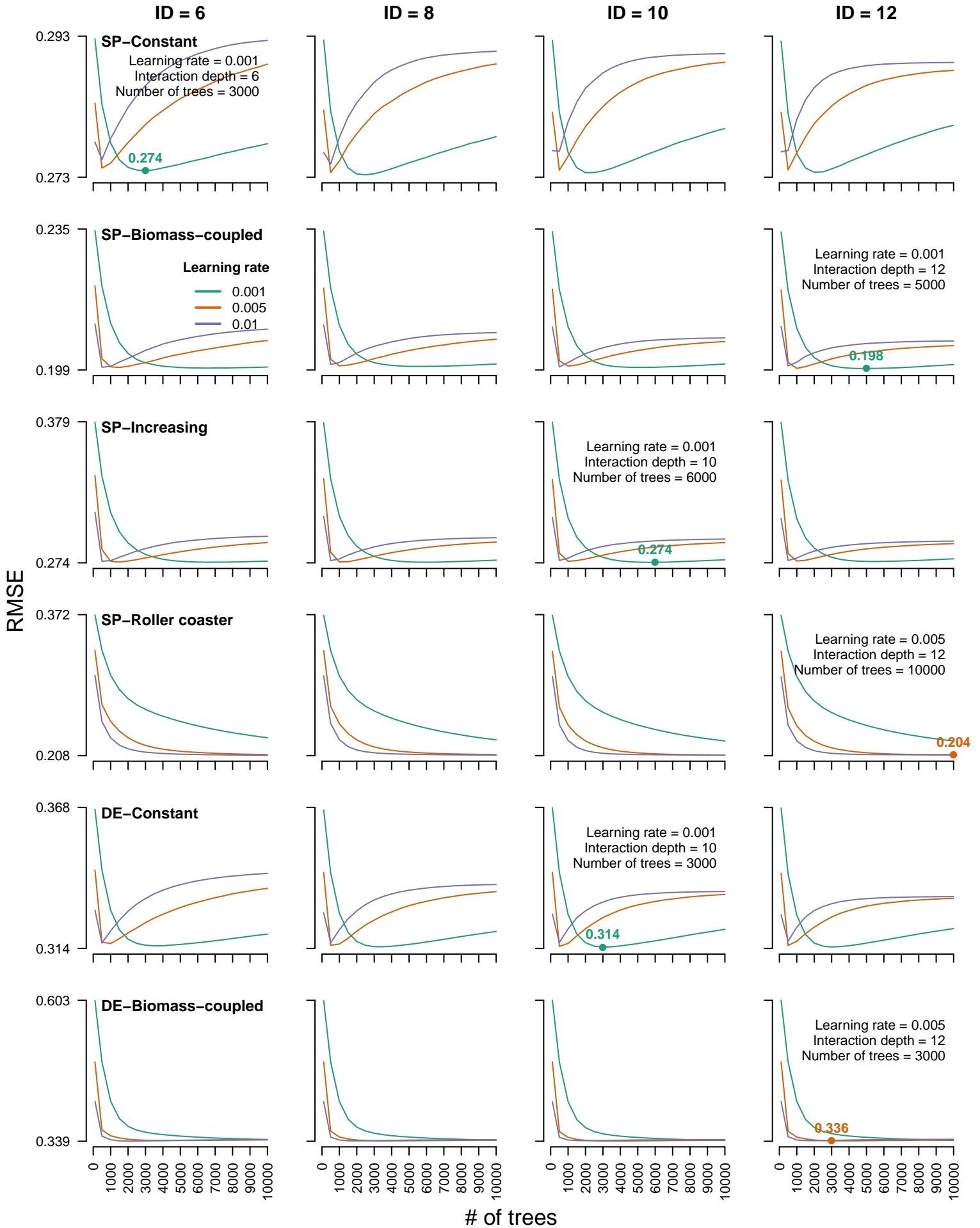
Supp. Figure 2. Model tuning curves showing the average root mean square error (RMSE) for each combination of candidate BRT model parameters (learning rate, interaction depth, # of trees) each of the four ED-tailored superensemble models. The optimal combination of model parameters (marked and labeled) is the combination that minimizes the RMSE.

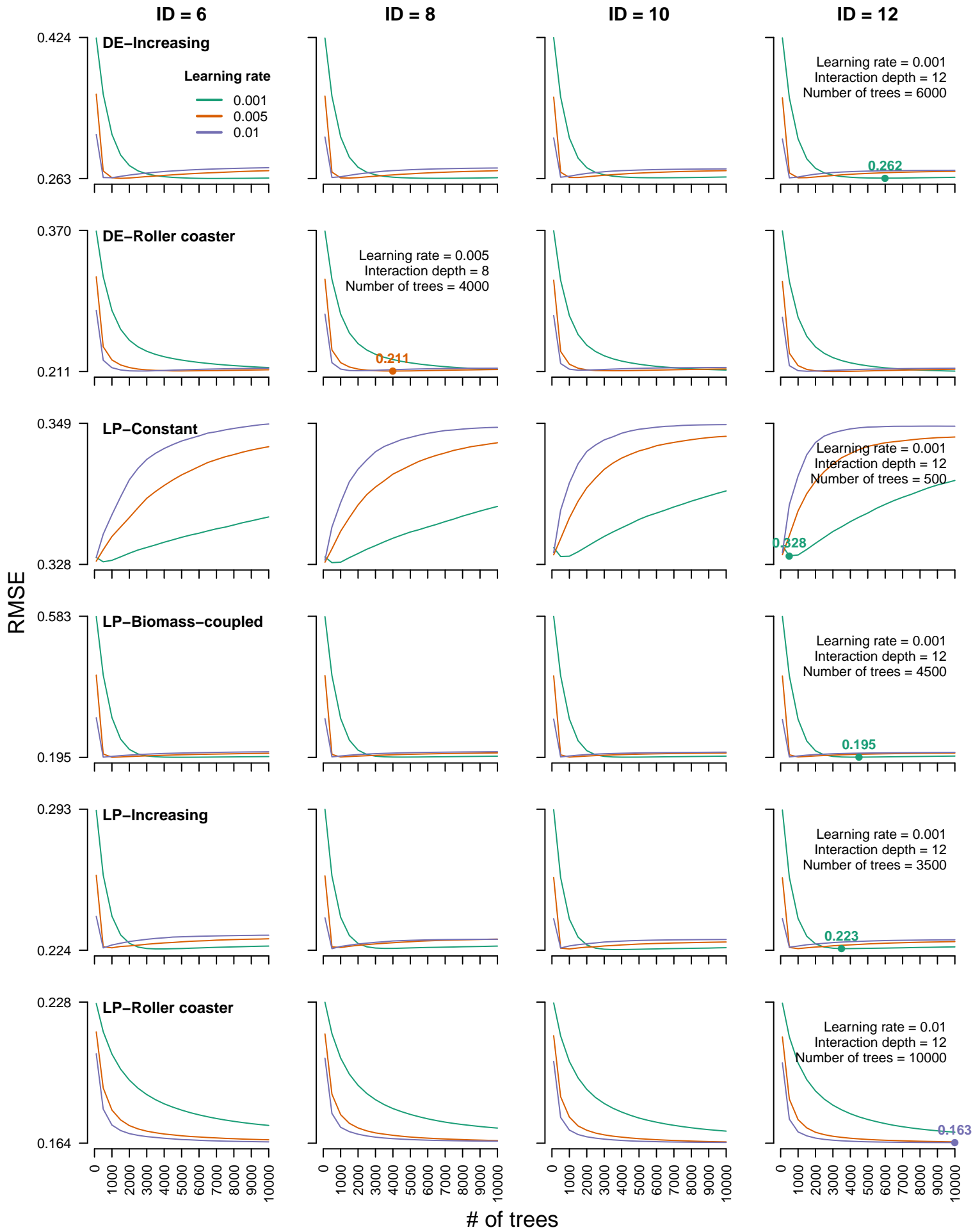
Appendix captions

Appendix A. Model tuning curves showing the average root mean square error (RMSE) for each combination of candidate BRT model parameters (learning rate, interaction depth, # of trees) for each of the twelve LH-ED-tailored superensemble models. The optimal combination of model parameters (marked and labeled) is the combination that minimizes the RMSE.

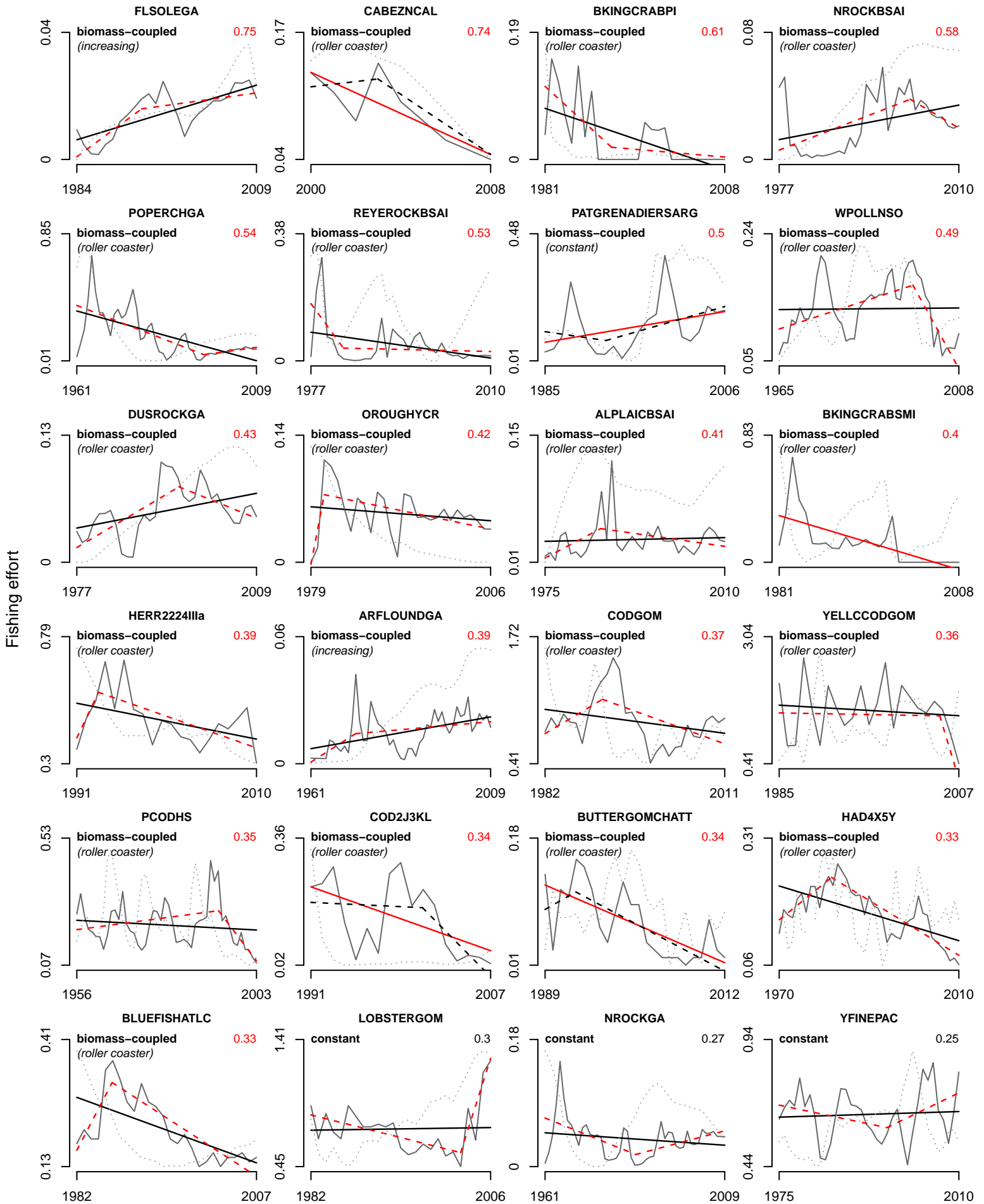
Appendix B. Illustration of the methods used to classify the effort dynamics experienced by the RAMLDB stocks. In each plot, the solid dark grey line indicates fishing effort over time and the dotted dark grey line indicates biomass over time. Early in the analysis, a Spearman's correlation between fishing effort and biomass lagged by one year greater than 0.3 was used to classify stocks experiencing "biomass-coupled" effort. The Spearman's correlations are printed in the top-right corner and are printed in red if greater than the 0.3 threshold. However, we decided that experts would be unable to accurately identify "biomass-coupled" effort and reclassified these stocks (the reclassification is listed in parentheses and italics) using linear and two-slope segmented regression. Linear and two-slope segmented regressions fit to the effort time series are shown as thick solid and thick dashed lines, respectively. The best regression model, determined through AIC, is shown in red. Details on how these regression models were used to classify effort dynamics are described in the text. The stocks are grouped by effort dynamics classification and sorted by descending effort-biomass correlation.

Appendix A.

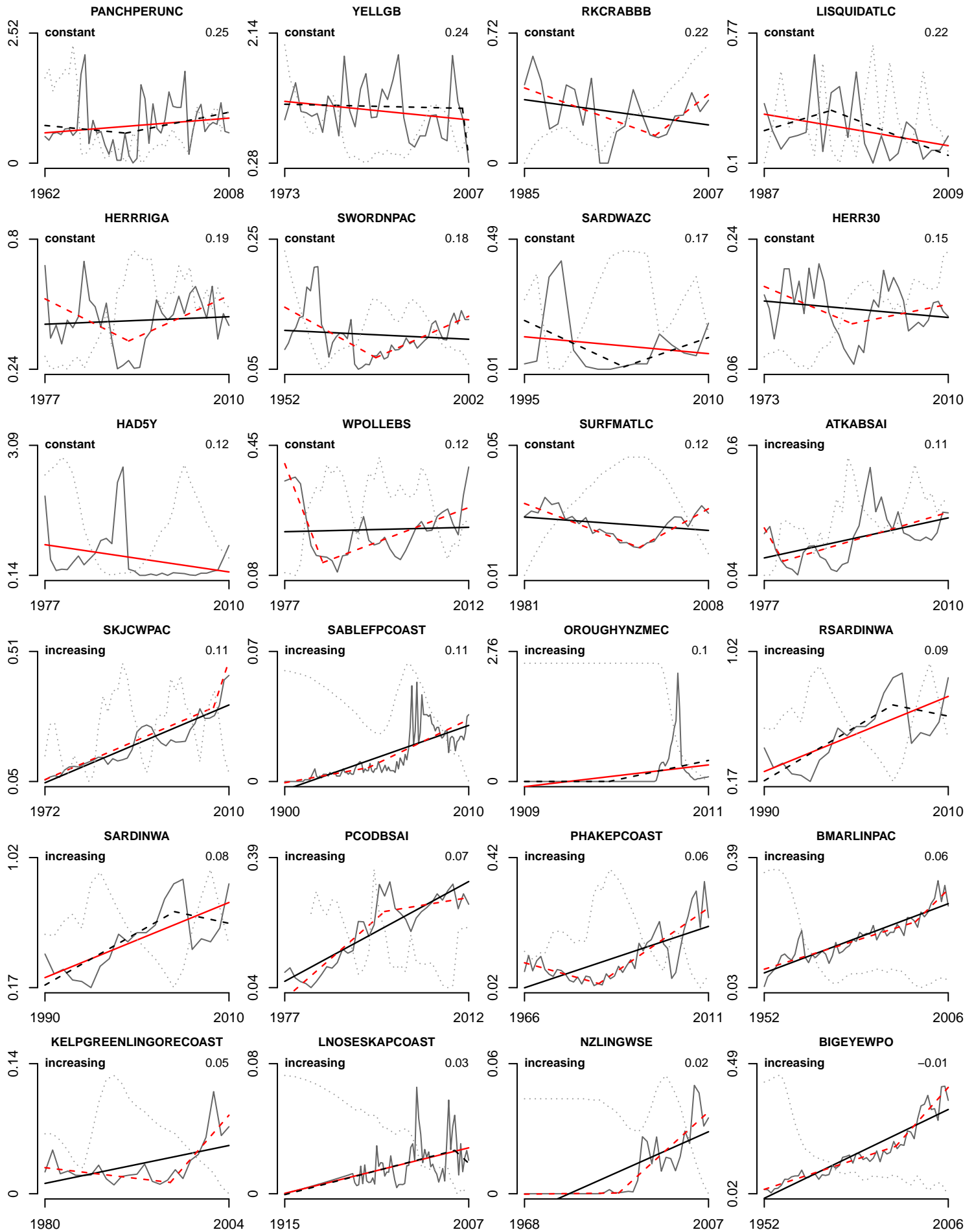




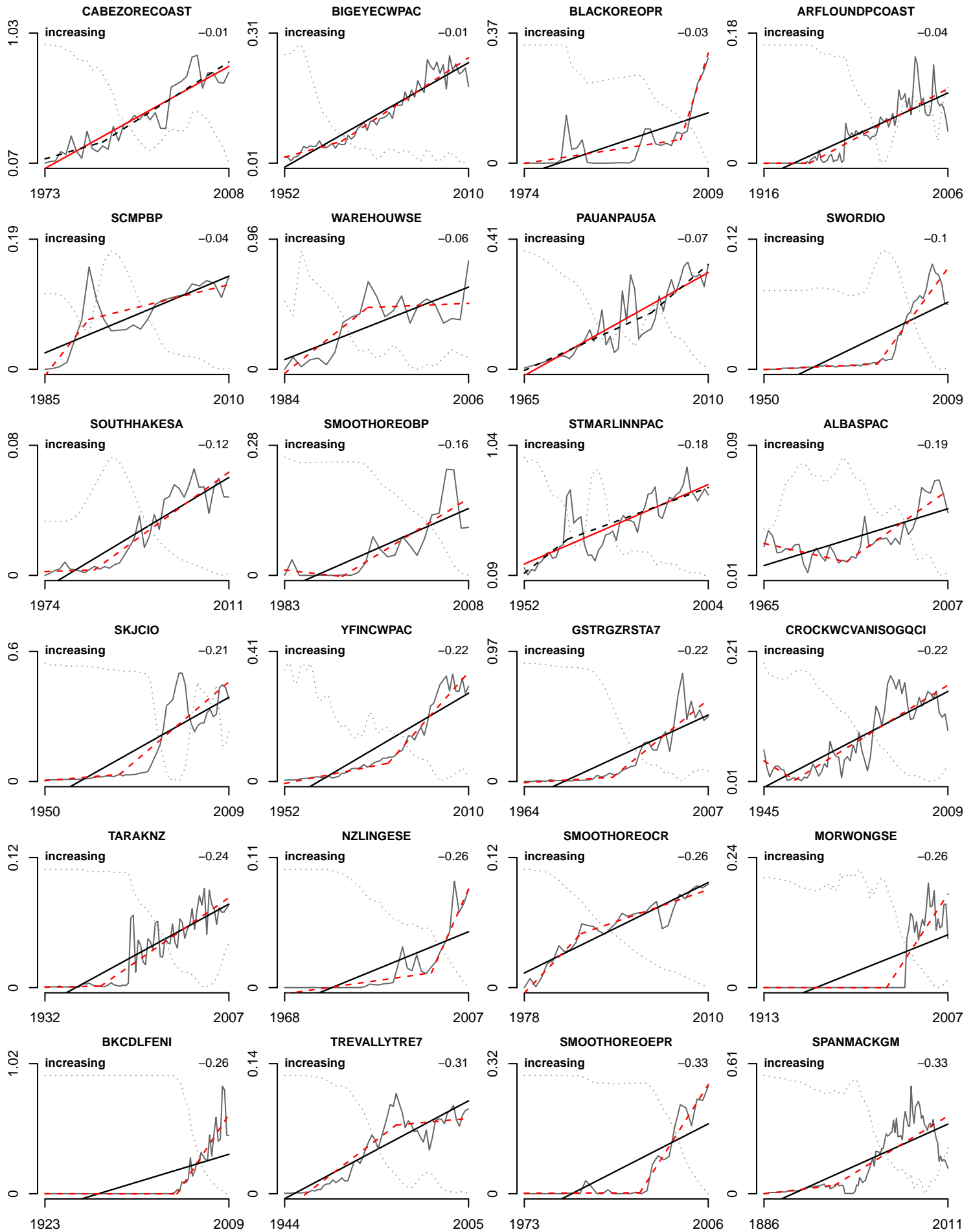
Appendix B.

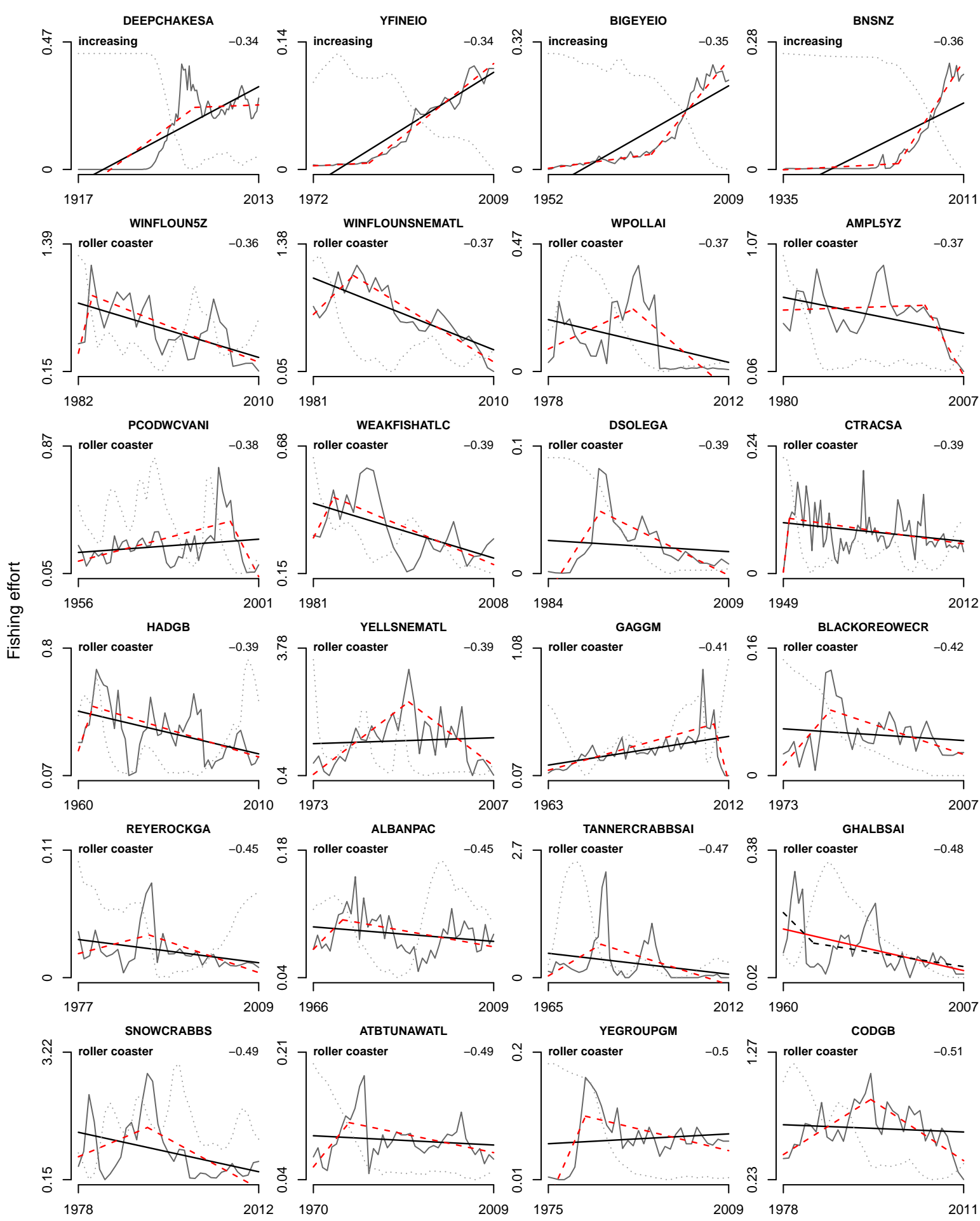


Fishing effort

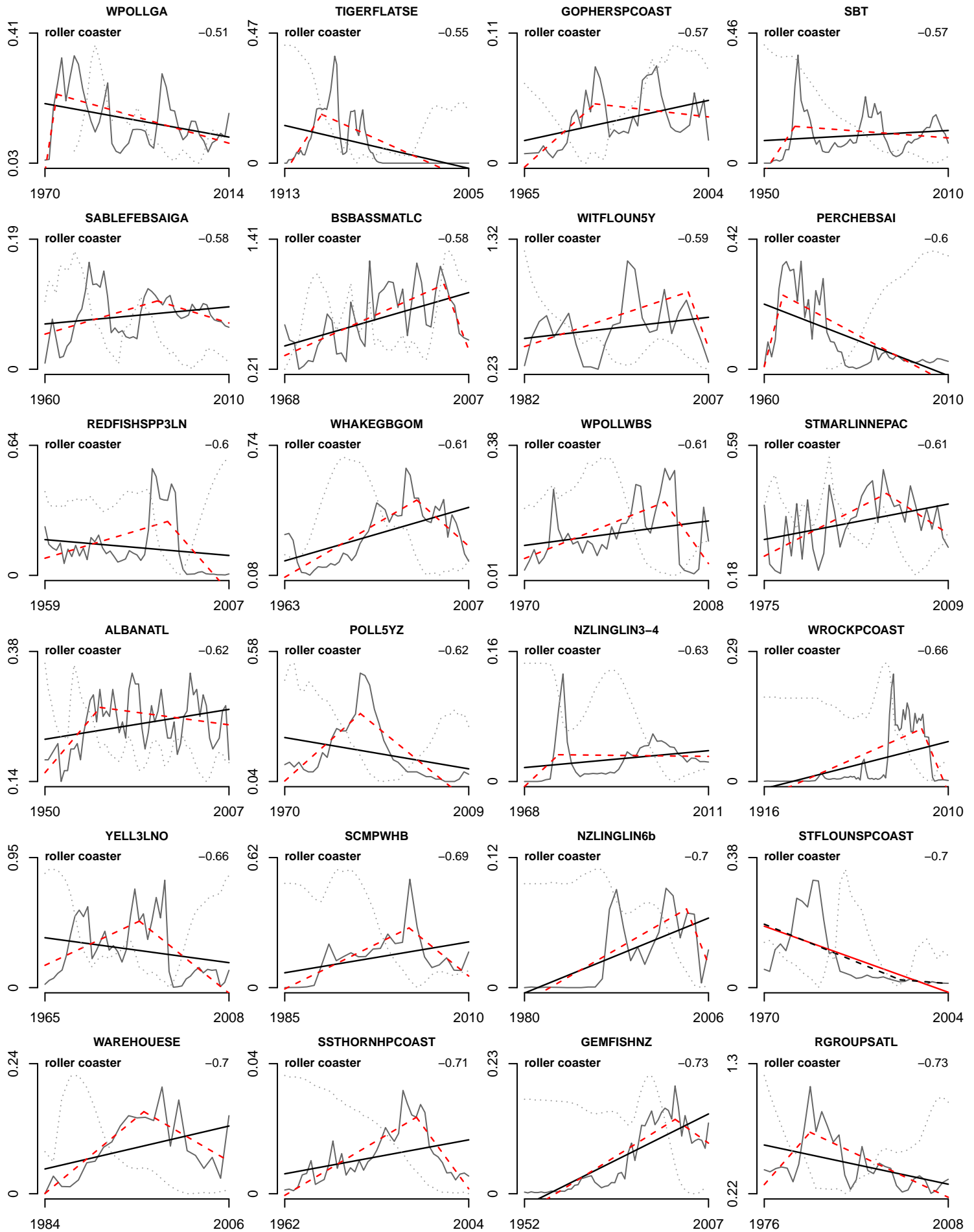


Fishing effort

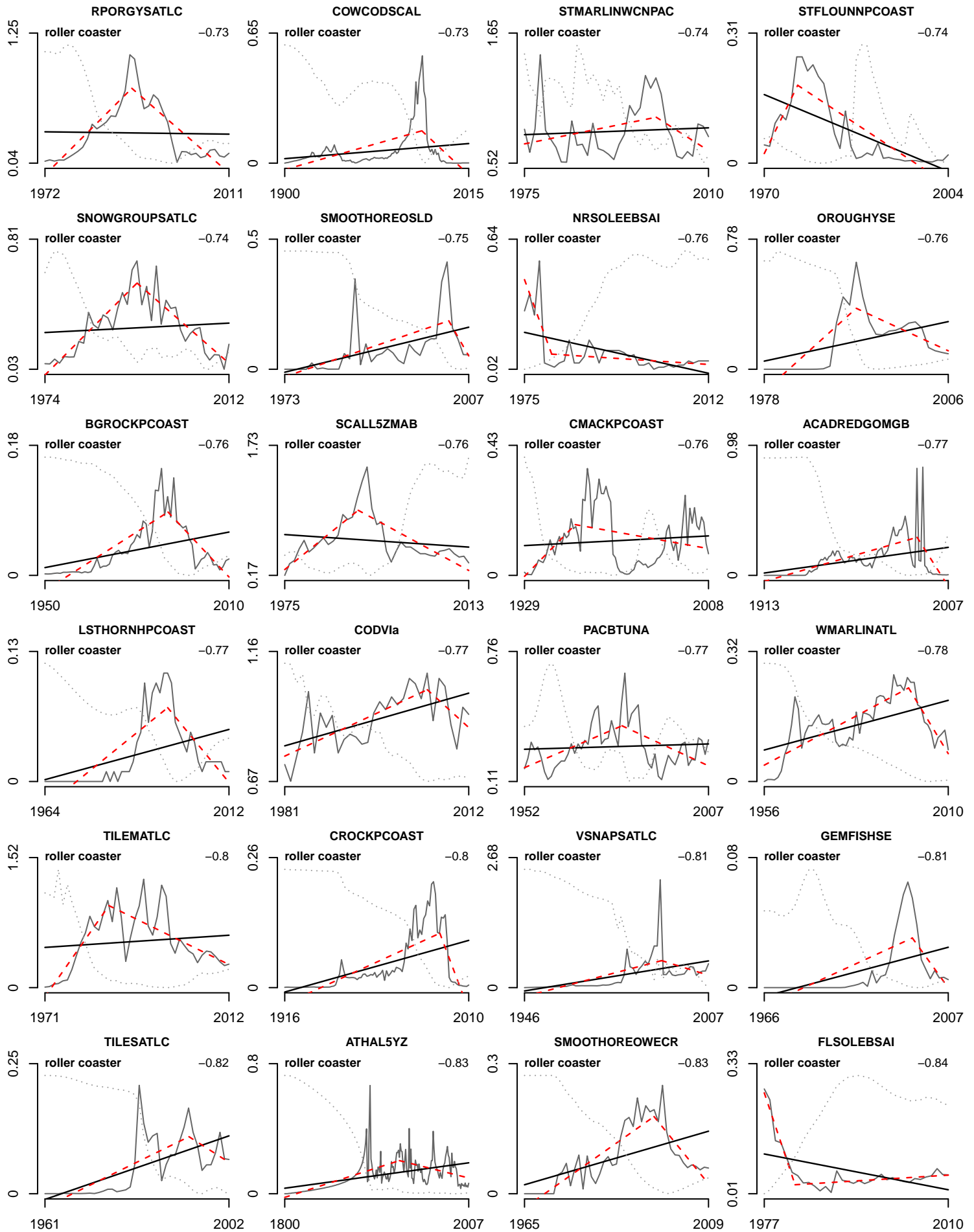




Fishing effort



Fishing effort



Fishing effort

