



ELSEVIER

Contents lists available at ScienceDirect

Journal of Marine Systems

journal homepage: www.elsevier.com/locate/jmarsys

Summary diagrams for coupled hydrodynamic-ecosystem model skill assessment

Jason K. Jolliff ^{a,*}, John C. Kindle ^b, Igor Shulman ^b, Bradley Penta ^b, Marjorie A.M. Friedrichs ^c, Robert Helber ^b, Robert A. Arnone ^b

^a Building 1009, Naval Research Laboratory, Stennis Space Center, (NRL-Stennis) Mississippi 39529, USA

^b NRL-Stennis, USA

^c Virginia Institute of Marine Science, P.O. Box 1346, Gloucester Point, VA 23062-1346, USA

ARTICLE INFO

Article history:
Received 30 April 2007
Accepted 2 May 2008
Available online xxx

Keywords:
Modeling
Marine ecosystem model
Statistical analysis
Remote sensing
Phytoplankton

ABSTRACT

The increasing complexity of coupled hydrodynamic-ecosystem models may require skill assessment methods that both quantify various aspects of model performance and visually summarize these aspects within compact diagrams. Hence summary diagrams, such as the Taylor diagram [Taylor, 2001, *Journal of Geophysical Research*, 106, D7, 7183–7192], may meet this requirement by exploiting mathematical relationships between widely known statistical quantities in order to succinctly display a suite of model skill metrics in a single plot. In this paper, sensitivity results from a coupled model are compared with Sea-viewing Wide Field-of-view Sensor (SeaWiFS) satellite ocean color data in order to assess the utility of the Taylor diagram and to develop a set of alternatives. Summary diagrams are only effective as skill assessment tools insofar as the statistical quantities they communicate adequately capture differentiable aspects of model performance. Here we demonstrate how the linear correlation coefficients and variance comparisons (pattern statistics) that constitute a Taylor diagram may fail to identify other potentially important aspects of coupled model performance, even if these quantities appear close to their ideal values. An additional skill assessment tool, the target diagram, is developed in order to provide summary information about how the pattern statistics and the bias (difference of mean values) each contribute to the magnitude of the total Root-Mean-Square Difference (RMSD). In addition, a potential inconsistency in the use of RMSD statistics as skill metrics for overall model and observation agreement is identified: underestimates of the observed field's variance are rewarded when the linear correlation scores are less than unity. An alternative skill score and skill score-based summary diagram is presented.

Published by Elsevier B.V.

1. Introduction

In general, mechanistic models that seek to simulate some natural phenomena must invariably be compared to observations in order to assess the model's skill. In accordance with this special volume on model skill assessment, we define *skill*

as the model's fidelity to the truth. We further presume that since the truth cannot be known, assessment of model skill must begin with a quantification of the misfit between model results and imperfect observations. An overview of various model skill metrics, which may include known statistical quantities or novel functions and mathematical techniques, is given in Stow et al. (submitted for publication). In this paper, we present a pragmatic evaluation of some widely known statistical quantities for the purpose of model skill assessment as well as how relationships between these quantities may be exploited to make compact diagrams that summarize multiple aspects of model performance, i.e., summary diagrams. An important component of this analysis is the relationship

* Corresponding author. Tel.: +1 228 688 5308; fax: +1 228 688 4149.
E-mail addresses: jolliff@nrlssc.navy.mil (J.K. Jolliff), kindle@nrlssc.navy.mil (J.C. Kindle), igor.shulman@nrlssc.navy.mil (I. Shulman), penta@nrlssc.navy.mil (B. Penta), marjy@vims.edu (M.A.M. Friedrichs), helber@nrlssc.navy.mil (R. Helber), bob.arnone@nrlssc.navy.mil (R.A. Arnone).

65 between various statistical quantities, which may be utilized
66 to produce summary diagrams, but may also be deceptive if
67 additional information is not presented. It is the general aim
68 of this paper to demonstrate that a comprehensive and bal-
69 anced approach to quantitative model skill assessment
70 should include, at the very least, an acknowledgement of
71 these relationships and an understanding of how they may
72 influence the appearance of model skill.

73 More specifically, however, summary diagrams may be
74 particularly suited to the task of skill assessment for spatially
75 complex models with multiple state variables, such as a
76 marine ecosystem model coupled to a hydrodynamic model
77 (coupled models – e.g., Franks and Chen, 2001; Gregg et al.,
78 2003; Walsh et al., 2003; Holt et al., 2005; Kindle et al., 2005;
79 Allen et al., 2007). Indeed, summary diagrams present a useful
80 method to succinctly communicate various aspects of coupled
81 model performance since extensive lists of metric values in
82 tabular form may become tedious. In addition, the use of
83 summary diagrams should also be encouraged in order to
84 address several other practical and scientific concerns. First,
85 many coupled model skill assessment exercises that have
86 appeared in recent literature still rely principally upon
87 graphics that emphasize the direct visual comparisons
88 between model results and observations (Stow et al., sub-
89 mitted for publication), such as a time series plot or a side-by-
90 side comparison of one to two-dimensional property fields
91 (chlorophyll, nitrate, etc.). If the statistical and graphical
92 techniques that are integral to the summary diagram approach
93 become more widely accepted and presented, then this may
94 encourage more quantitative statements about coupled model
95 skill. Second, summary diagrams are particularly useful for
96 quantitatively comparing the performance of an ensemble of
97 different models or multiple permutations of a single model.
98 Given that there remains continuing uncertainty in the struc-
99 ture and parameterization of ecosystem models (e.g., Frie-
100 drichs et al., 2007), summary and quantitative skill assessment
101 techniques may become an efficient facilitator of improved
102 prognostic performance.

103 Accordingly, one potential statistical and graphical skill
104 assessment approach is to render a Taylor diagram (Taylor,
105 2001). Taylor diagrams exploit relationships between known
106 statistical quantities in order to provide summary information
107 about particular aspects of model performance and were
108 developed to aid in the monitoring of complex ocean–atmo-
109 sphere climate models. The Taylor diagram, as is the case for
110 many potential model skill assessment tools, is not discipline
111 specific, and several recent marine ecosystem modeling papers
112 have presented them as part of a model skill assessment
113 scheme (Gruber et al., 2006; Raick et al., 2007). Here we begin
114 with an assessment of the Taylor diagram and the statistics it
115 communicates for the specific purpose of coupled model skill
116 assessment. Taylor diagrams are an appropriate place to begin
117 our evaluation of summary diagrams given their increasing use
118 in a wide range of modeling disciplines; however, summary
119 diagrams are only as useful as the metrics they communicate,
120 and so our analysis includes an exposition of how relationships
121 between widely known statistical quantities may be further
122 utilized to construct other types of summary diagrams that
123 communicate additional aspects of model performance.

124 While the statistical methods and diagrams developed and
125 discussed here may potentially be applied to many other

126 types of model result to data comparisons, we nonetheless
127 present results from a coupled hydrodynamic–ecosystem
128 model and ocean color products derived from SeaWiFS sate-
129 llite ocean color data in order to explicitly illustrate potential
130 problems arising from this type of skill assessment. To that
131 end, summary information about the modeling and satellite
132 ocean color methods is given below (Section 2), whereas
133 detailed description of statistical methods and display
134 techniques are fully explicated in due course of the main
135 analysis (Section 3). In Section 3.1, we examine the Taylor
136 diagram and the univariate statistics it summarizes by pre-
137 senting several example applications that demonstrate the
138 strengths and weaknesses of this approach. In Section 3.2, we
139 develop an alternative summary diagram, the target diagram,
140 which provides information about additional aspects of
141 model performance that may be of particular concern to the
142 skill assessment of ecosystem models. In Section 3.3, we
143 identify a potentially undesirable property of RMSD-based
144 metrics, and present an alternative skill score and skill score-
145 based summary diagram.

2. Methods 146

147 Results from an experimental ecosystem modeling environ-
148 ment, the Naval Research Laboratory Ecological-Photoche-
149 mical-Bio-Optical-Numerical Experiment (which for brevity
150 is referred to as Neptune), are presented here as a prototypical
151 example of a complex modeling system. Detailed description
152 of the Neptune modeling construct, including all state equa-
153 tions, parameter designations, and optical calculations, may
154 be found in Jolliff and Kindle (2007). The modeling system is
155 composed of four core elements: (1) the biogeochemical
156 model that describes the flow and transformation of ele-
157 mental reservoirs (carbon, nitrogen, and phosphorus) as a
158 result of phytoplankton primary production and subsequent
159 physiological processes and trophic interactions; (2) a visible
160 optics module that relates the biogeochemical elemental
161 reservoirs to spectrally explicit optical properties, describes
162 the vertically resolved attenuation of incident, spectrally de-
163 composed irradiance, and budgets photons absorbed by living
164 phytoplankton to perform light-growth calculations; (3) an
165 ultraviolet (UV) optics module that determines the attenua-
166 tion of spectrally decomposed UV irradiance and the potential
167 UV-stimulated photochemical degradation of colored dis-
168 solved organic matter (CDOM); and (4) a description of the
169 spectrally decomposed UV and visible irradiance boundary
170 conditions.

171 The Neptune system is designed for integration with any
172 hydrodynamic model capable of describing the advection–
173 diffusion of state variables. Here we examine the one-dimen-
174 sional case by coupling the model to the Modular Ocean Data
175 Assimilation System (MODAS). MODAS is described in Fox et al.
176 (2002). Briefly, the system uses optimal interpolation (Breth-
177 erton et al., 1976) to render daily satellite estimates of sea
178 surface temperature (SST) and sea surface height (SSH) onto a
179 two-dimensional grid. A subsurface temperature profile is then
180 retrieved from the U.S. Navy's Master Oceanographic Observa-
181 tional Data Set. Deviation from subsurface climatology is then
182 estimated based upon SST and SSH deviation from surface
183 climatology. The result is a synthetic three-dimensional tem-
184 perature field.

185 The MODAS fields were averaged over 4 years (2001–
186 2004) to approximate an average annual cycle of summer
187 thermal stratification followed by winter overturn for a $1^\circ \times 1^\circ$
188 area in the western Gulf of Mexico (center position 24.0° N,
189 94.5° W). Vertical eddy diffusion coefficients were imputed
190 from MODAS synthetic temperature fields using the Paca-
191 nowski and Philander (1981) vertical mixing scheme. Daily
192 and vertically resolved (total depth (z) = 161 m; $\Delta z = 1$ m) eddy
193 diffusion coefficients were used to solve for the vertical tur-
194 bulent mixing of model state variables using a fully implicit
195 method with a time step of 1800 s. The coupled model was
196 initialized using temperature–nutrient relationships ob-
197 served in the Gulf of Mexico (Jochens et al., 2002) and then
198 run for ten simulation years to solve for the steady state
199 solution for transformations of carbon, nitrogen, and phos-
200 phorus in the upper ocean. The system was forced to material
201 conservation by implicit remineralization of all particulates
202 that sank below the deepest grid cell ($z = 161$ m).

203 The coupled model results were compared to local area
204 coverage SeaWiFS ocean color data that were received and
205 archived at the Naval Research Laboratory (NRL), Stennis Space
206 Center. The satellite data were processed and the intervening
207 atmospheric signal removed using NRL's Automated Processing
208 System (APS). The atmospheric correction procedures are com-
209 pliant with National Aeronautics and Space Administration
210 SeaWiFS data processing protocols. Three NRL APS products
211 derived from SeaWiFS data were examined: (1) the surface
212 chlorophyll-*a* concentration, which was determined from the
213 OC4v4 band ratio algorithm (O'Reilly et al., 1998); (2) the surface
214 phytoplankton absorption coefficient (443 nm); and (3) the
215 surface colored detrital matter (CDM) absorption coefficient
216 (412 nm). The latter two products were determined from the
217 multiband quasi-analytic algorithm (Lee et al., 2002), which
218 estimates total absorption coefficients over SeaWiFS visible
219 bands and then further decomposes them into phytoplankton
220 and detrital contributions. Each daily spatial mean of SeaWiFS
221 data through 4 years (2001–2004) from the 1° western Gulf of
222 Mexico grid was used to construct a satellite ocean color time
223 series wherein missing days due to clouds were accounted for

224 via linear interpolation. The time series was lowpass filtered to
225 remove variability from frequencies higher than 10 days; the
226 averages were then computed to construct the annual
227 climatology.

228 3. Results

229 The model results are compared with the daily climatol-
230 ogy calculated from 4 years of SeaWiFS data (Fig. 1) for three
231 surface bio-optical fields: the surface chlorophyll-*a* concen-
232 tration, the surface phytoplankton absorption coefficient
233 (443 nm), and the surface CDM absorption coefficient
234 (412 nm). The satellite estimate of these surface quantities
235 will be herein referred to as the reference field and the
236 model's simulated surface bio-optical quantities will be
237 referred to as simply the model field.

238 The Neptune model's three size-based phytoplankton
239 functional groups are presently parameterized so that pico-
240 phytoplankton have a higher absorption efficiency (per unit
241 chlorophyll-*a*) than larger phytoplankton, as has been observed
242 in the laboratory and in the field (e.g., Bricaud et al., 2004;
243 Millan-Nunez et al., 2004). Thus the model phytoplankton
244 absorption and total chlorophyll fields may vary with respect to
245 one another due to differences in the relative dominance of
246 simulated phytoplankton size fractions. In the example given in
247 the following section, the satellite estimates of phytoplankton
248 absorption and chlorophyll are thus used as a potential ob-
249 servational constraint on the simulated competition between
250 phytoplankton size fractions.

251 3.1. Taylor diagrams and pattern statistics

252 For the one-dimensional case wherein the model's surface
253 values are averaged over the upper 10 m each simulated day
254 and are compared with a single daily reference value, the
255 model and reference fields resemble sinusoidal functions of
256 time, or waveforms (Fig. 1). Analogously, a measure of the
257 potential phase shift between the two waveforms is also more
258 generally a common measure of the agreement between two

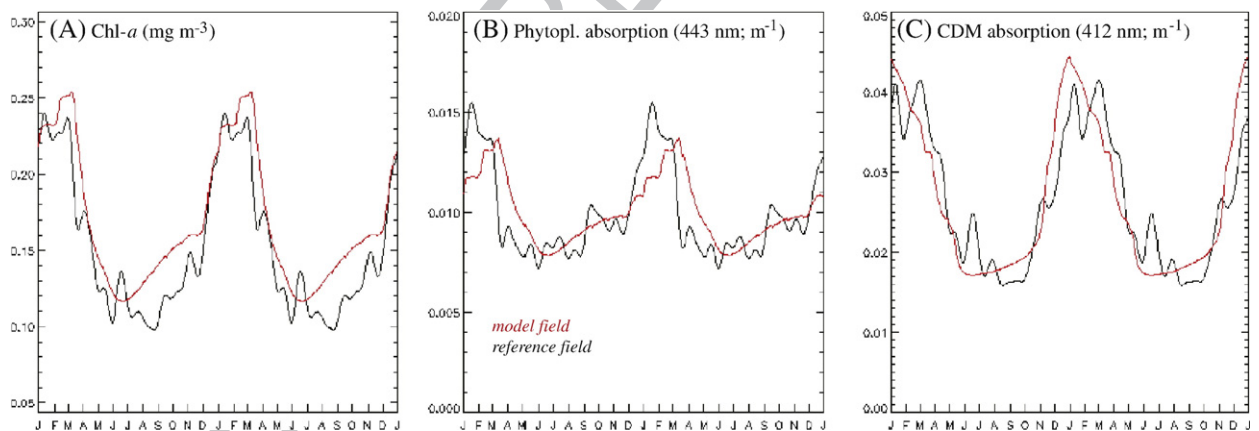


Fig. 1. Daily surface values for the (A) chlorophyll-*a* concentration (mg m^{-3}), (B) phytoplankton absorption coefficient (443 nm, m^{-1}), and (C) CDM absorption coefficient (412 nm, m^{-1}) are indicated for the final 2 years of the model's steady state solution (red line) and the SeaWiFS climatology (black line). Two years are shown in order to emphasize the winter peak and bring further emphasis to temporal misfits (i.e., phase misfits quantified by linear correlation coefficients). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

259 fields: the linear correlation coefficient, R , which is defined
260 by:

$$R = \frac{\frac{1}{N} \sum_{n=1}^N (m_n - \bar{m})(r_n - \bar{r})}{\sigma_m \sigma_r} \quad (1)$$

262 The letter m indicates the model field, r indicates the re-
263 ference field, the overbar indicates the average, and σ is the
264 standard deviation.

265 The correlation coefficient is bounded by the range
266 $-1.0 \leq R \leq 1.0$. In general, as the phase between two temporal
267 signals approaches agreement, R approaches 1.0. It is difficult,
268 however, to discern information about the differences in
269 amplitude between two signals from R alone. For this reason,
270 another summary statistic, the normalized standard deviation,
271 may be introduced:

$$\sigma^* = \frac{\sigma_m}{\sigma_r} \quad (2)$$

274 The normalized standard deviation and the correlation
275 coefficient from each of the three model to reference field
276 comparisons may be displayed on a single Taylor diagram
277 (Fig. 2). The Taylor diagram is a polar coordinate diagram that
278 assigns the angular position to the inverse cosine of the correla-
279 tion coefficient, R . A correlation coefficient of 0 is thus 90°
280 away from a correlation coefficient of 1 (see scaling on Fig. 2).
281 The radial (along-axis) distance from the origin is assigned to
282 the normalized standard deviation, σ^* . The reference field
283 point, which is comprised of the statistics generated from a
284 redundant reference to reference comparison, is indicated for
285 the polar coordinates (1.0, 0.0). The model to reference com-

286 parison points may then be gauged by how close they fall to the
287 reference point. This distance is proportional to the unbiased
288 Root-Mean-Square Difference (RMSD'), as defined by:

$$\text{RMSD}' = \left(\frac{1}{N} \sum_{n=1}^N [(m_n - \bar{m}) - (r_n - \bar{r})]^2 \right)^{0.5} \quad (3)$$

290 where the overbars indicate the mean. The term *unbiased* is
291 used herein to emphasize that Eq. (3) removes any information
292 about the potential bias (B), which is defined as the difference
293 between the means of the two fields:

$$B = \bar{m} - \bar{r} \quad (4)$$

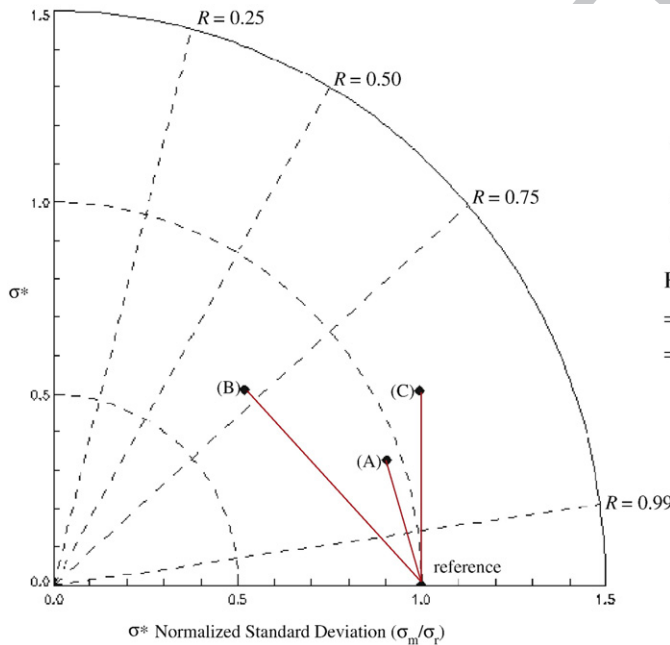
296 In other words, the unbiased RMSD (RMSD') is equal to the
297 total RMSD if there is no bias between the model and
298 reference fields. This may be verified given the quadratic
299 relationship between the unbiased RMSD, the bias, and the
300 total RMSD:

$$\text{RMSD}^2 = B^2 + \text{RMSD}'^2 \quad (5)$$

303 where the total RMSD is a measure of the average magnitude
of difference and is defined by:

$$\text{RMSD} = \left[\frac{1}{N} \sum_{n=1}^N (m_n - r_n)^2 \right]^{0.5} \quad (6)$$

306 In contrast, the unbiased RMSD may be conceptualized as
307 an overall measure of the agreement between the amplitude
308 (σ) and phase (R) of two temporal patterns. For this reason,
309 the correlation coefficient (R), normalized standard deviation
310 (σ^*), and unbiased RMSD are collectively referred to herein as
311



$$B^* = \frac{(\bar{m} - \bar{r})}{\sigma_r} = 0.385 \text{ (A)}$$

$$B^* = 0.062 \text{ (B)}$$

$$B^* = 0.037 \text{ (C)}$$

$$\text{RMSD}^{*'} = \sqrt{1.0 + \sigma^{*2} - 2\sigma^*R} = 0.340 \text{ (A)}$$

$$= 0.701 \text{ (B)}$$

$$= 0.510 \text{ (C)}$$

Fig. 2. Taylor diagram rendering of the model to reference field comparisons shown in Fig. 1: (A) chlorophyll- a concentration (mg m^{-3}), (B) phytoplankton absorption coefficient (443 nm, m^{-1}), and (C) CDM absorption coefficient (412 nm, m^{-1}). As explained in the text, the radial distance is proportional to the normalized standard deviation (σ^*) and the angular position corresponds to the linear correlation coefficient (R values). In accordance with Eq. (7), the distances between the labeled points and the reference point are proportional to the unbiased RMSD, Eq. (3).

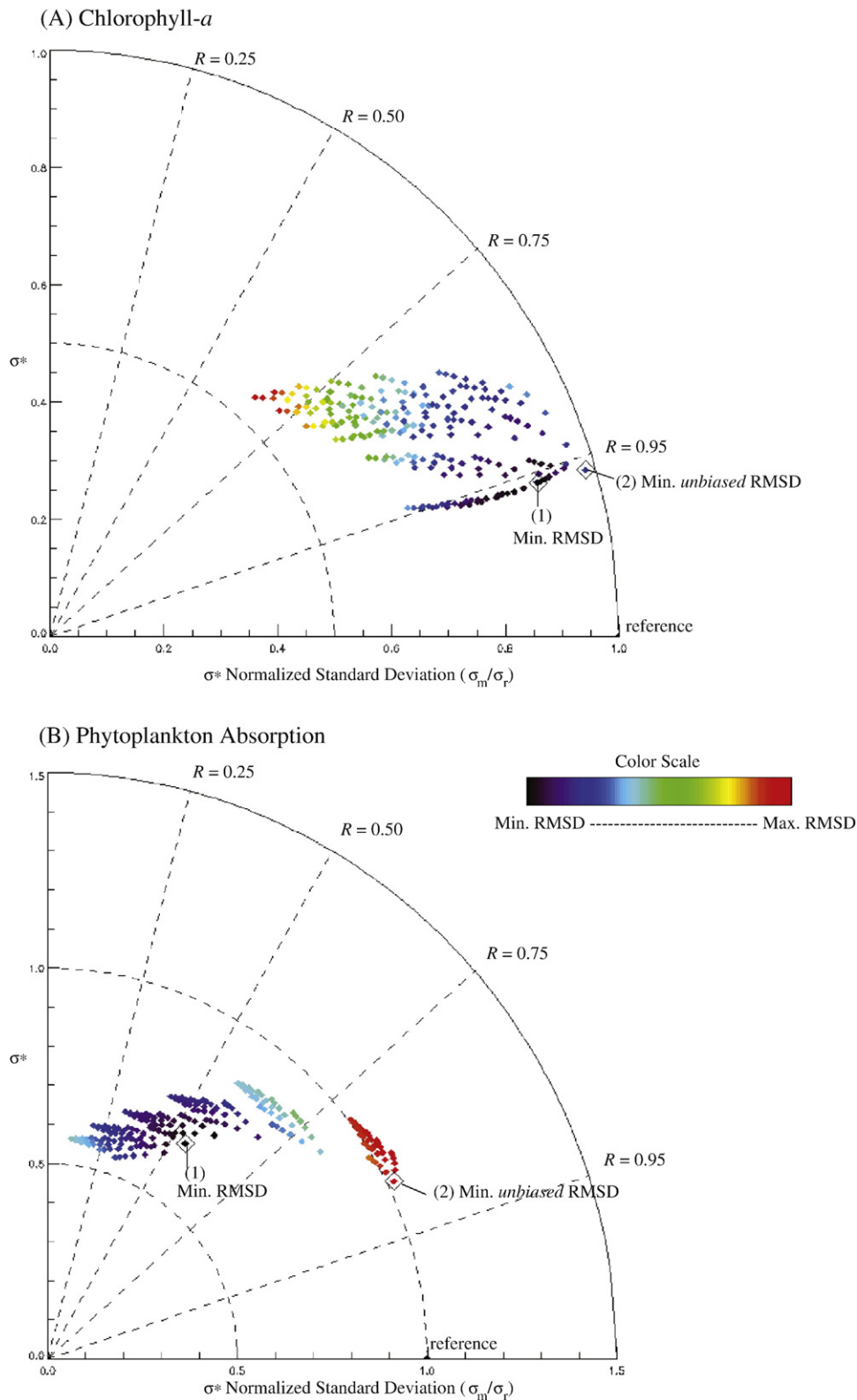


Fig. 3. Taylor diagrams for grazing sensitivity model executions showing model to reference statistics for the (A) surface chlorophyll-*a* field and (B) the surface phytoplankton absorption field. The minimum total RMSD (1) and the minimum unbiased RMSD (2) are indicated on each plot. The color scale is added to both Taylor diagrams and corresponds to the minimum total RMSD (black) to the maximum total RMSD (red) for each set of model to reference comparison statistics. The time series results corresponding to points (1) and (2) in (B) are shown in Fig. 4. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

335 pattern statistics. The three pattern statistics are related to one
336 another by:

$$\text{RMSD}^2 = \sigma_r^2 + \sigma_m^2 - 2\sigma_r\sigma_m R \quad (7)$$

338 It is this relationship that makes the Taylor diagram useful:
339 the individual contribution of misfits in amplitude may be
340 compared to misfits in phase to discern how they contribute
341 to the unbiased RMSD. Since the diagram is in standard
342 deviation normalized space, the distance from the model
343 points to the reference points is also proportional to Eq. (7),
344 which recast in standard deviation normalized units (indi-
345 cated by the asterisk) becomes:

$$\text{RMSD}^{*f} = \sqrt{1.0 + \sigma^{*2} - 2\sigma^*R} \quad (8)$$

346 Note also that it can be shown that the minimum of this
348 function occurs where $\sigma^* = R$. This is an important relationship
349 that we will refer to at several points later in the text.

350 Fig. 2 shows that the chlorophyll model to reference field
351 comparison point (A) appears closest to the reference point,
352 whereas the phytoplankton absorption comparison point (B)
353 appears farthest due to a poorer correlation as well as an
354 underestimate of the standard deviation. Indeed, the chloro-
355 phyll comparison has the lowest normalized and unbiased
356 RMSD. However, the normalized bias, defined as:

$$B_* = \frac{(\bar{m} - \bar{r})}{\sigma_r} \quad (9)$$

358 is much larger for the model chlorophyll field, which con-
359 sistentlly tends to overestimate the reference field (as shown
360 in Fig. 1A). Thus caution must be applied when interpreting a
361 Taylor diagram wherein no information about the bias is
362 included.

The importance of adding information about the bias may
also be further demonstrated using a large number of model
executions, such as during a sensitivity analysis. The advan-
tage of the Taylor diagram in such cases is that it allows one to
discern how the phase and amplitude of a simulated field
change as the model is modified. The disadvantage is that
information about any potential model to reference field bias
must be somehow added to the diagram.

For example, the mortality rate for phytoplankton (ε_r) in
the Neptune ecological model is described using the Ivlev
(1961) formulation:

$$\varepsilon_r = \varepsilon_m \left(1.0 - e^{-Iv(C)}\right) \quad (10)$$

where Iv is the Ivlev parameter that describes how the maxi-
mum potential mortality rate (ε_m) is attenuated with decreasing
phytoplankton biomass (C). With three phytoplankton func-
tional groups and an estimated Iv parameter space incremented
for 6 values, there are 216 potential grazing permutations.

The results of 216 separate model executions are shown on
two Taylor diagrams (Fig. 3). For brevity, only the first two field
comparisons, phytoplankton chlorophyll and phytoplankton
absorption, are shown since the CDM absorption field is
somewhat less sensitive to the grazing parameter selections. It
is important to note that the model and reference fields were
not log-transformed. In this case, it would not make a con-
siderable difference; however, if there were large outliers in
either field then log-transformation may significantly impact
the value of statistical quantities. Some investigators may
choose to log-transform the fields first, particularly if the bio-
optical fields range over several orders of magnitude. If the
fields are log-transformed then the investigator should be
cognizant that statistical quantities generated from non log-
transformed values may be different.

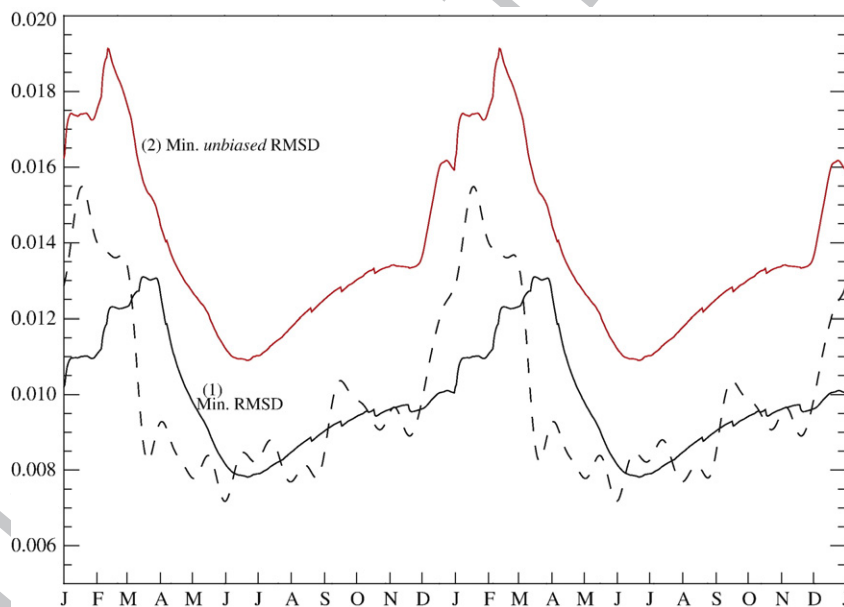


Fig. 4. The reference field phytoplankton absorption (dashed line) is compared to the minimum total RMSD (1 – solid black line) and the minimum unbiased RMSD (2 – red line); these time series correspond to points (1) and (2) in Fig. 3B. As in Fig. 1, two years are shown to emphasize the winter peak and draw emphasis to phase misfits quantified by the linear correlation coefficients. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

414 In both Taylor diagrams presented here, the model points
 396 that come closest to the reference point have the smallest
 397 unbiased RMSD value (Fig. 3). It would appear that the cluster
 398 of model points closest to the reference point may thus
 399 provide the closest fit to the data. Here, however, the
 400 inclusion of a relative total RMSD color scale, which indicates
 401 the range of minimum to maximum total RMSD using a
 402 spectral (rainbow) color scaling increment (Fig. 3), reveals
 403 that some points nearest the reference point may have larger
 404 total RMSD values. This is particularly the case for phyto-
 405 plankton absorption (Fig. 3B) where the cluster of points
 406 closest to the reference point also have the largest total RMSD
 407 values. For the phytoplankton absorption field, improvement
 408 in the correlation coefficient appears to come at the expense
 409 of an increase in the bias, and consequently, the total RMSD.
 410 The minimum total RMSD (point 1) and minimum unbiased
 411 RMSD (point 2) from the phytoplankton absorption compar-
 412 isons are also shown as a time series plot (Fig. 4). Clearly, the
 413 red line (minimum unbiased RMSD) has a better phase agree-
 414 ment but overestimates the observed values.

415 In coupled hydrodynamic-ecosystem modeling applica-
 416 tions, information about the bias and the total RMSD may be
 417 just as important to the investigator as information about the

418 pattern statistics, particularly when evaluating the sensitivity
 419 of a model to parameter selection for the purpose of mini-
 420 mizing the magnitude of the misfit between the model and
 421 reference fields. Taylor (2001) suggested adding lines of
 422 various lengths corresponding to the total RMSD in propor-
 423 tion to the unbiased RMSD onto the Taylor diagram; however,
 424 this procedure may result in a confusing diagram when large
 425 numbers of model runs are compared. A color scale modi-
 426 fication of the Taylor diagram, as shown here (Fig. 3), may also
 427 be useful but the overall import of the Taylor diagram may
 428 nonetheless be easily misinterpreted.

3.2. Target diagrams

429
 430 An alternative to the Taylor diagram is to formulate a
 431 target diagram that provides summary information about the
 432 pattern statistics as well as the bias thus yielding a broader
 433 overview of their respective contributions to the total RMSD.
 434 The relationship between the bias, unbiased RMSD, and the
 435 total RMSD (Eq. (5)) provides a convenient starting point to
 436 construct such a diagram. In a simple Cartesian coordinate
 437 system, the unbiased RMSD may serve as the X-axis and the
 438 bias may serve as the Y-axis. The distance between the origin

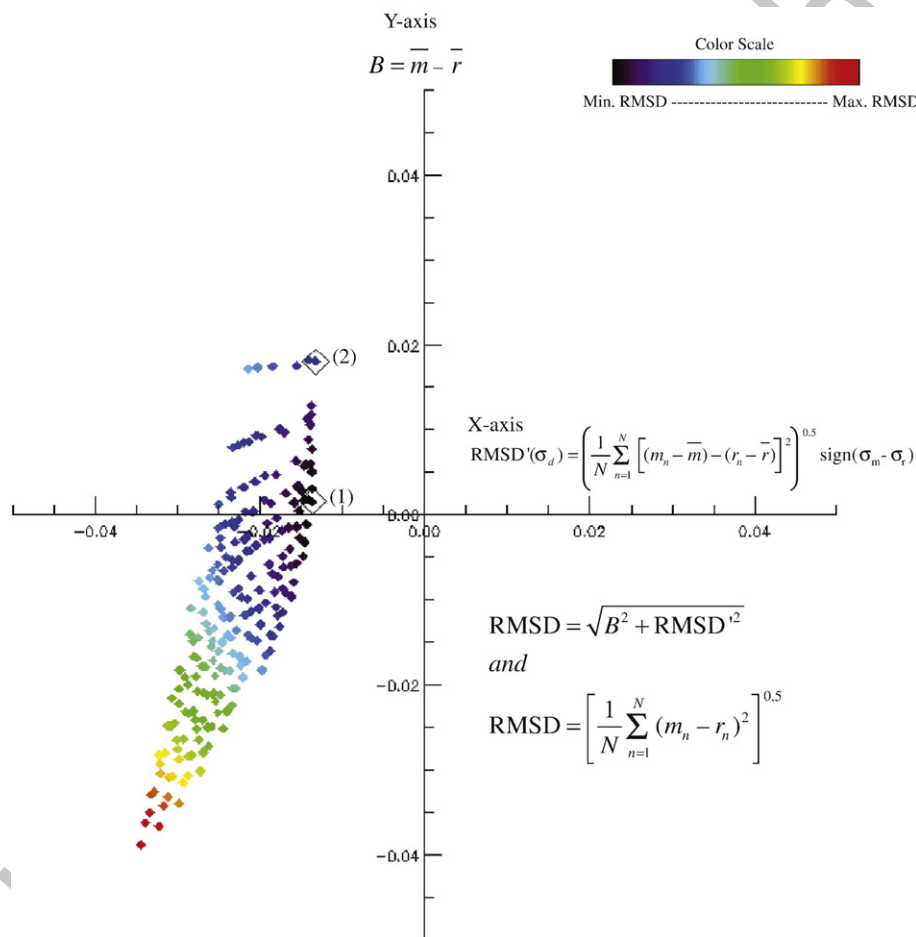


Fig. 5. Target diagram for model chlorophyll-*a* and reference chlorophyll-*a* comparisons. The Y-axis corresponds to the bias, the X-axis corresponds to the unbiased RMSD multiplied by the sign of the model and reference standard deviation difference (σ_d), and the distance from each point to the origin is proportional to the total RMSD. The minimum total RMSD (1) and the minimum unbiased RMSD (2) are indicated on the plot. The color scaling is the same as in Fig. 3.

439 and the model versus observation statistics (any point, s ,
440 within the X, Y Cartesian space) is then equal to the total RMSD
441 (Fig. 5).

442 By definition, the X -axis (unbiased RMSD) must always
443 be positive. However, the $X < 0.0$ region of the Cartesian
444 coordinate space may be utilized if the unbiased RMSD is
445 multiplied by the sign of the standard deviation difference
(σ_d):

$$\sigma_d = \text{sign}(\sigma_m - \sigma_r) \quad (11)$$

449 The resulting target diagram thus provides information
450 about whether the model standard deviation is larger
451 ($X > 0$) or smaller ($X < 0$) than the reference field's standard
452 deviation, in addition to a positive ($Y > 0$) or negative bias
453 ($Y < 0$) (Fig. 5). The units of this diagram are all in chlo-
454 rophyll concentration (mg m^{-3}), but this may again be
455 addressed by normalizing the quantities by the reference

456 field standard deviation (Fig. 6), such that the distance of
457 each point from the origin is the standard deviation nor-
458 malized total RMSD:

$$\text{RMSD}^{*2} = B^{*2} + \text{RMSD}^{*r2} \quad (12)$$

460 Rendering the diagram in normalized units allows one to
461 better compare the model's chlorophyll performance with
462 other potential areas of performance such as CDM absorp-
463 tion and phytoplankton absorption.

464 Furthermore, markers within the diagram may be added to
465 provide an additional basis for interpreting model perfor-
466 mance. For example, the investigator may wish to gauge how
467 the model's total RMSD compares to the time series mean. In
468 other words, if the first guess is the time series average, does
469 the model provide an overall improvement over the first guess
470 with respect to the minimization of the average misfit bet-
471 ween the model and reference fields?

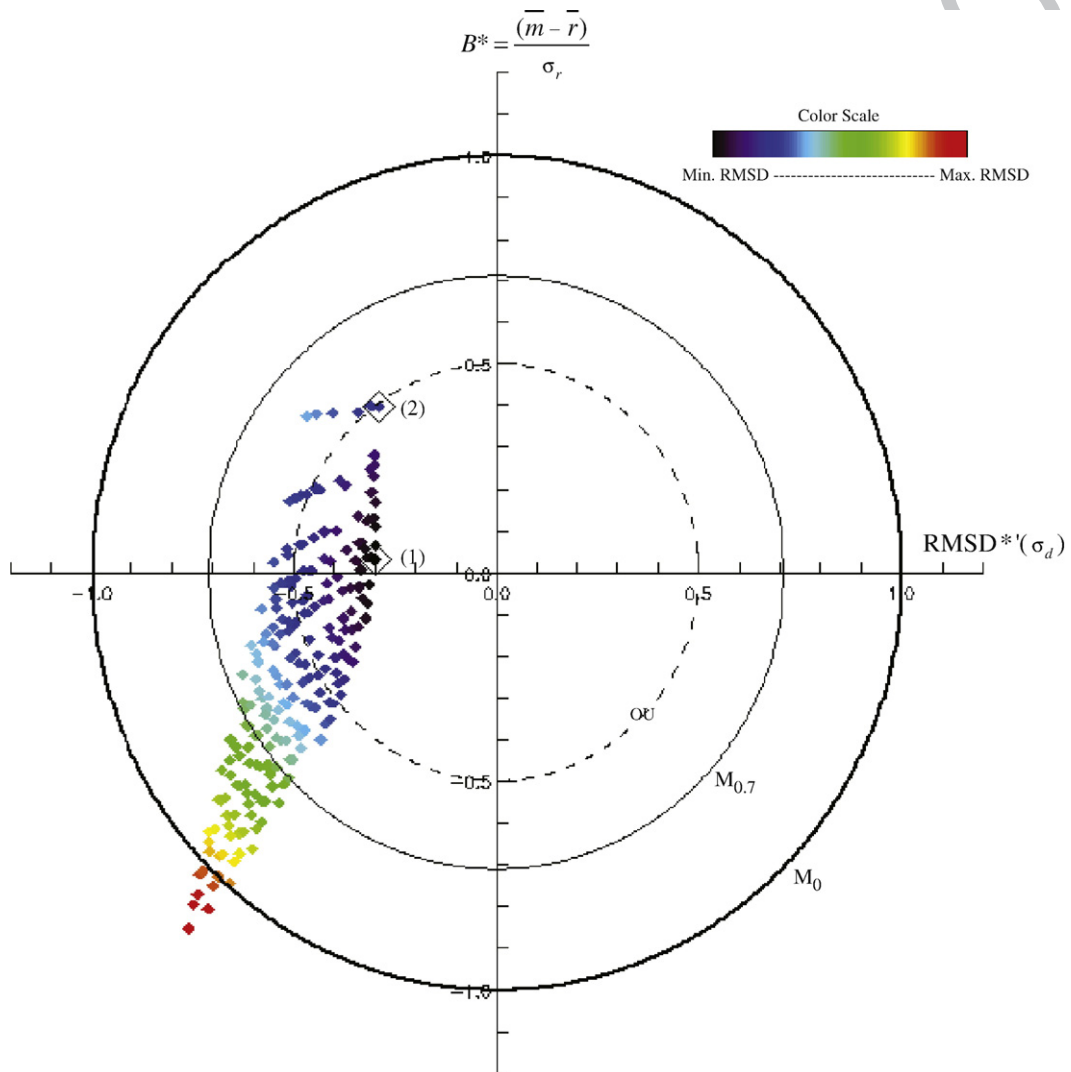


Fig. 6. Normalized target diagram for model chlorophyll- a and reference chlorophyll- a comparisons. The axes are the same as in Fig. 4, only they are normalized by the reference field standard deviation (indicated by *). The thick line (M_0) corresponds to a normalized total RMSD of 1.0, the thin line ($M_{0.7}$) corresponds to $\text{RMSD}^* = 0.71$. The significance of these markers is explained in the text. The dashed line represents the threshold of observational uncertainty (OU). The minimum total RMSD (1) and the minimum unbiased RMSD (2) are indicated on the plot. The color scaling is the same as in Figs. 3 and 5.

422 The total RMSD between the reference field and the
423 reference field mean is simply the reference field's standard
424 deviation. Since the diagram is in standard deviation normal-
425 ized units, a normalized total RMSD value of 1.0 provides a
426 convenient performance marker (marker M_0 , Fig. 6). If the
427 investigator is concerned with the total RMSD, and not merely
428 the pattern statistics, then any points greater than $\text{RMSD}^* = 1$
429 may be considered poor performers since they offer no im-
430 provement over the time series average.

431 It is also interesting to note that the normalized total RMSD
432 (RMSD^*) is related to the modeling efficiency (MEF) metric
433 presented in Stow et al. (submitted for publication) via the
434 relationship: $\text{MEF} = 1 - \text{RMSD}^{*2}$. The MEF may be used to discern
435 how well a model performs as a predictor of the data compared
436 to the mean of the data (Stow et al., 2003; Nash and Sutcliffe,
437 1970). This underscores the significance of the $\text{RMSD}^* = 1$ (M_0)
438 marker within the normalized target diagram since points be-
439 tween it and the origin also have a better than average MEF score.

440 A weakness of the target diagram is that it does not provide
441 explicit information about the correlation coefficient. However,
442 there are certain limits inherent in the statistics summarized by
443 the diagram that one may use to make some inference about
444 the correlation coefficient. For example, recall the relationship
445 between the correlation coefficient, the normalized standard
446 deviation, and the normalized and unbiased RMSD (Eq. (8)). It
447 can be shown that for values of R (where $-1.0 \leq R < 0.0$) the
448 minimum value of RMSD^* for all potential values of σ^* (where
449 $0.0 < \sigma^* < \infty$) approaches 1.0. Thus no model/reference compar-
450 ison points that appear on the target diagram within the range of
451 $-1.0 < X < 1.0$ can be negatively correlated. Since the square of the
452 normalized bias must always be positive, then by extension all
453 points where $\text{RMSD}^* < 1.0$ must also be positively correlated.
454 In other words, the first marker at $\text{RMSD}^* = 1.0$ (marker M_0 , Fig. 5)
455 also establishes that all points between it and the origin are
456 positively correlated. Positively correlated results may appear
457 outside this marker; however, these points will have a large
458 magnitude of difference from the observations due to either a
459 significant bias, a difference in variance, or some combination
460 thereof. This relationship may be formally expressed as follows:

$$\text{for } \forall s \in \{\text{RMSD}^* \mid \text{RMSD}^* \leq 1.0\} \rightarrow R > 0.0 \quad (13)$$

512 where s is a notation for any point on the target diagram. Similar
513 such markers based upon the correlation coefficient may be
514 established closer to the origin for values of R where $R > 0.0$. In
515 accordance with Eq. (8), the minimum value of RMSD^* occurs
516 for any positive value of R where $\sigma^* = R$. Thus if one wants to
517 determine the minimum unbiased RMSD value possible (M_{R1})
518 given a specific correlation value, $R1$, then the solution may be
519 expressed as:

$$M_{R1} = \min(\text{RMSD}^*) = \sqrt{1.0 + R1^2 - 2R1^2} \quad (14)$$

521 Since the minimum total RMSD must also occur where the
522 bias is equal to 0.0, M_{R1} is also the minimum total RMSD
523 value for a given correlation coefficient value, $R1$. For the
524 general case where $R1 > 0.0$:

$$\text{for } \forall s \in \{\text{RMSD}^* \mid \text{RMSD}^* \leq M_{R1}\} \rightarrow R \geq R1 \quad (15)$$

527 For example, Fig. 6 shows the second marker towards the
528 origin for $R1 = 0.7$. Thus all points between this marker ($M_{0.7}$)

and the origin are indicative of a correlation coefficient
greater than 0.7.

531 The color scale in Fig. 6 is redundant: both the distance
532 from the origin and the color index are proportional to the
533 total RMSD. The color variable is thus left as a free variable
534 that may be used to also explicitly indicate the correlation
535 coefficient, or it may be used to indicate any supplemental
536 information regarding the simulations that are displayed in
537 the diagram (Friedrichs et al., submitted for publication). In
538 our example, the sensitivity analysis is focused upon the
539 grazing parameters. We may define an aggregate index of
540 phytoplankton grazing stress (AI) as the sum of the three Ivlev
541 parameters and display this index using the color scale, as in
542 Fig. 7. Clearly, the AI most appreciably impacts the bias: as
543 aggregate grazing stress increases the simulations consis-
544 tently underestimate the satellite-based observations of
545 surface chlorophyll. Furthermore, the lowest aggregate graz-
546 ing stress corresponds to the highest bias (point 2, Fig. 7).

547 Diagrams that summarize repeated comparisons of model
548 results and data should also make some indication of un-
549 certainties that exist within the data. One may define data as
550 truth plus some unknown observational uncertainty. The ad-
551 vantage of using a satellite climatology based upon a large
552 number of spatial means, as in this case, is that one may
553 choose to assume that the ensemble average observational
554 uncertainty approaches zero as the total number of observa-
555 tions becomes very large ($\sim n > 1000$). One approach might
556 be to state that assumption and forego any further indication
557 of observational uncertainty. A note of caution must also be
558 applied insofar as this approach assumes that the observa-
559 tional uncertainty is also unbiased.

560 Nevertheless, for the more general case there exists a large
561 sum of potential observational uncertainties arising, in part,
562 from measurement error. For satellite data, these errors may
563 arise from imperfections in the satellite sensor, errors in the
564 algorithms applied, atmospheric correction errors, and nume-
565 rous other areas beyond the scope of this paper. It is therefore
566 reasonable to assume that there must be some average mini-
567 mum threshold value for the total RMSD below which further
568 improvement in model/data agreement may not be signifi-
569 cant. The dashed line in Fig. 5 is an estimate of this observa-
570 tional uncertainty (OU) threshold. Points that fall between
571 this limit and the origin are all within the range of estimated
572 observational uncertainty.

573 To be sure, observational uncertainty is a potentially com-
574 plicated and contentious subject. Our objective here is to simply
575 represent some estimate of this uncertainty on the target
576 diagram so as to indicate where further efforts towards im-
577 proved model to data agreement may not be a prudent use of
578 time and resources. While it is entirely reasonable and appro-
579 priate to assume that observational uncertainty does provide an
580 upper-limit upon potential improvements in model perfor-
581 mance, our tentative estimates of this average uncertainty
582 should be regarded as preliminary and much more work in this
583 area needs to be done.

584 In this case, an average observational uncertainty was
585 assumed for the satellite time series based on literature values
586 for chlorophyll algorithm accuracy in optically deep waters
587 (Bailey and Werdell, 2006; McClain et al., 2006) without any
588 further consideration of the uncertainty within the measure-
589 ments to which the satellite data are compared. If the average

590 observational uncertainty (α) is expressed as a percent, then $\alpha\bar{r}$
 591 may be used as an estimate for the average value of uncertainty
 592 for the time series. For example, a α value of $\pm 15\%$ and an
 593 average chlorophyll-*a* observation of 0.2 mg m^{-3} would yield an
 594 average uncertainty of $\pm 0.03 \text{ mg m}^{-3}$. A model to reference field
 595 total RMSD of $< 0.03 \text{ mg m}^{-3}$ is within the average observational
 596 uncertainty threshold and further improvement (model to data
 597 misfit reduction) may not be meaningful.

598 This assumed OU limit may be placed on the target diagram
 599 by normalizing $\alpha\bar{r}$ by the reference field standard
 600 deviation (dashed line, Fig. 7). The normalization procedure
 601 effectively means that the assumption of average observa-
 602 tional uncertainty (α) is divided by the coefficient of variation,
 603 which is the reference field standard deviation divided by the
 604 reference field mean. The coefficient of variation is a common
 605 measure of the dispersion within a distribution. It is beyond
 606 the scope of this paper to further examine how the dispersion,

in turn, may be impacted by the observational uncertainty, 607
 but we recognize that they are not necessarily independent. 608

In summary, the target diagram displays the model to 609
 reference field bias (*Y*-axis) and the model to reference field 610
 unbiased RMSD (*X*-axis). The distance between any point, *s*, 611
 and the origin is then the value of the total RMSD. All of the 612
 quantities may be normalized by the reference field standard 613
 deviation to remove the units of measurement. The outermost 614
 marker ($M_0 = \text{RMSD}^* = 1.0$) establishes that all points between 615
 it and the origin represent positively correlated model and 616
 reference fields, and also have a better than average MEF score. 617
 A second marker may be added to indicate another positive *R* 618
 value, such as $R = 0.7$, for which all points between it and the 619
 origin are greater than *R*. Finally, a dashed line indicates the 620
 estimate of average observational uncertainty and further 621
 model to data agreement for points between this marker and 622
 the origin may not be meaningful. 623

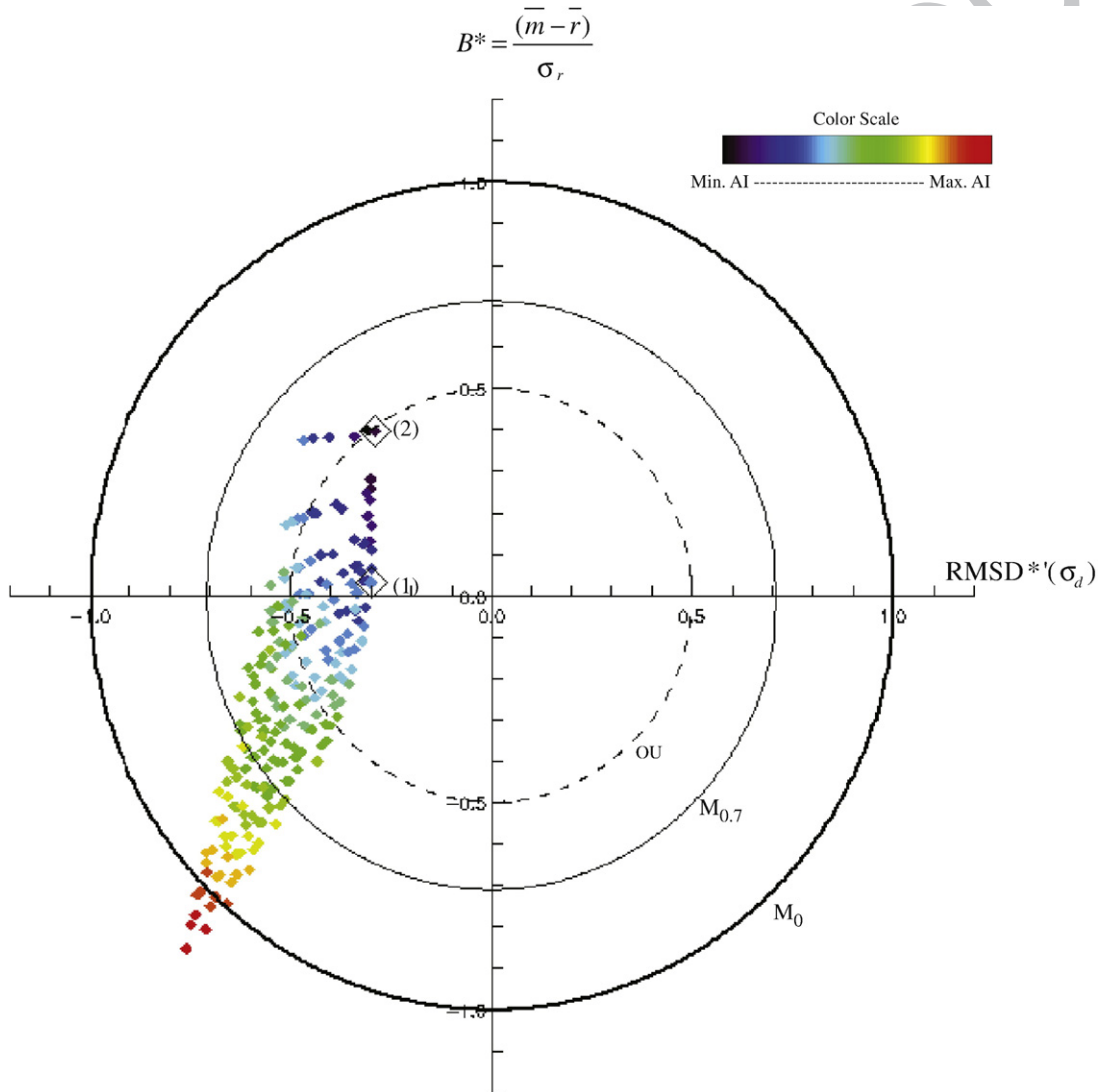


Fig. 7. Normalized target diagram for model chlorophyll-*a* and reference chlorophyll-*a* comparisons. The axes and the markers are the same as in Fig. 6. The color scaling has been changed to indicate the aggregate index (AI) for grazing stress, as explained in the text.

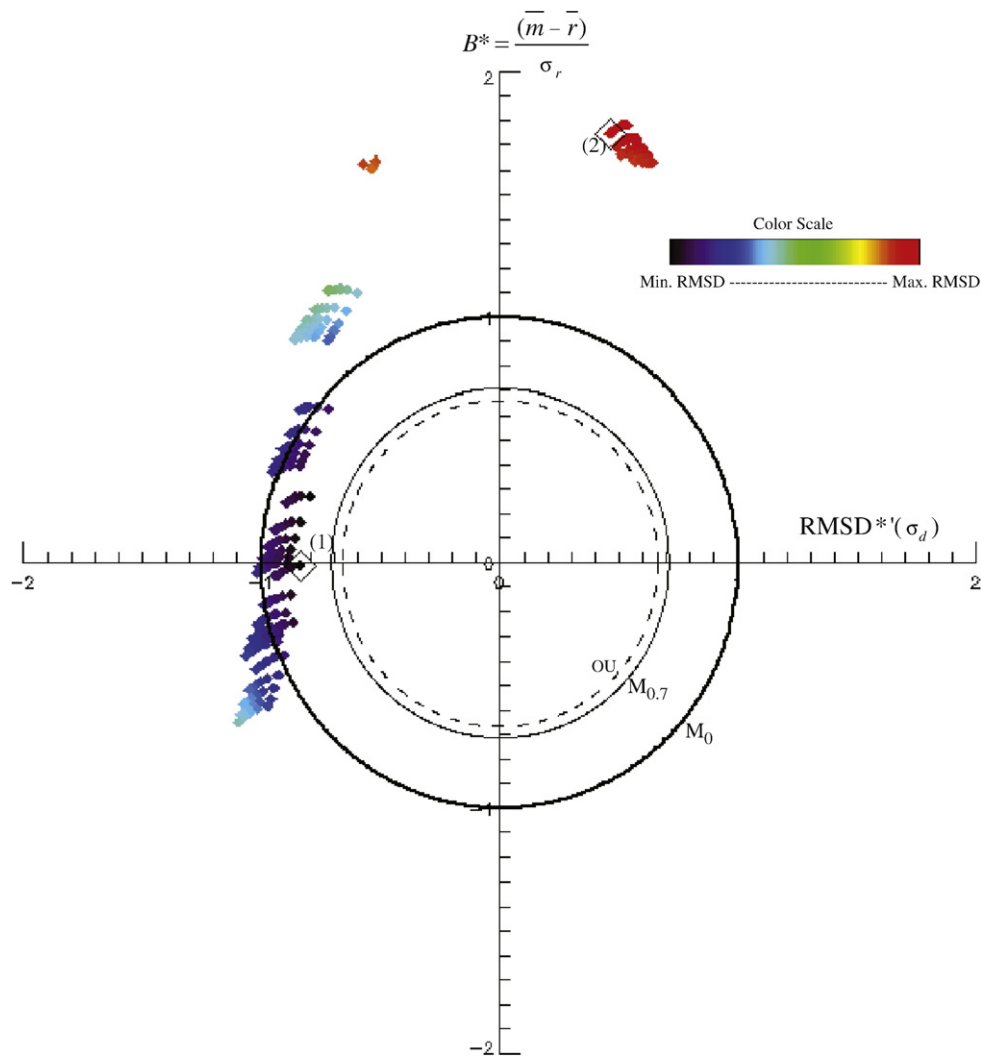


Fig. 8. Normalized target diagram for model/reference phytoplankton absorption fields. The axes are normalized by the reference field standard deviation (indicated by *). The thick line (M_0) corresponds to a normalized total RMSD of 1.0, the thin line ($M_{0.7}$) corresponds to $\text{RMSD}^* = 0.71$. The significance of these markers is explained in the text. The dashed line represents the threshold of observational uncertainty (OU). The minimum total RMSD (1) and the minimum unbiased RMSD (2) are indicated on the plot.

704 The target diagram was also constructed for the phyto- 641
 625 plankton absorption field (Fig. 8). In order to display the entire 642
 626 set of model versus reference comparisons for phytoplankton 643
 627 absorption, the scale for the target diagram (Fig. 8) had to be 644
 628 expanded to encompass $\text{RMSD}^* = 2$. Note that the simulations 645
 629 with the best pattern statistics (Fig. 3B) also have a very large 646
 630 positive bias (red cluster, Fig. 8). In this particular case, the 647
 631 target diagram better delineates poor performing model execu- 648
 632 tions than the Taylor diagram since the model is prone to a 649
 633 large bias for this field.

634 3.3. The skill target diagram

635 Additional alternatives to the Taylor diagram for summar- 650
 636 izing pattern statistics as a measure of model skill may be 651
 637 preferable since there is a subtle discrepancy between improv- 652
 638 ing the unbiased RMSD and improving the individual 653
 639 correlation coefficient and standard deviation statistics, and 654
 640 there may be circumstances where this consideration is im- 655

portant. For example, consider that there may be fundamental 641
 limits to the expected agreement between a model and a 642
 reference field. Even if all model inaccuracies and observa- 643
 tional uncertainties could be eliminated, there may yet remain 644
 unforced oscillations that prevent exact model/reference field 645
 agreement. Suppose that an estimate of this uncertainty yields 646
 a maximum potentially attainable correlation coefficient of 647
 0.65. As stated in Section 3.1, the minimum value of the un- 648
 biased RMSD occurs where $\sigma^* = R$ for positive values of R . 649

This relationship may be demonstrated on a Taylor diagram 650
 (Fig. 9). For $R = 0.65$ the minimum RMSD^* value occurs where 651
 $\sigma^* = 0.65$. The three sets of pattern statistics correspond to the 652
 waveforms in Fig. 9B. The minimum average difference is the 653
 smallest amplitude pattern, but if amplitude and phase are 654
 weighed equally, as in a potential alternative measures of model 655
 skill, then the waveform where $\sigma^* = 1$ may be the most skillful. 656

This example demonstrates the implicit contradiction between 657
 minimizing the RMSD and improving σ^* towards an ideal 658
 value of 1.0. If the goal is to improve the total RMSD then σ^* 659

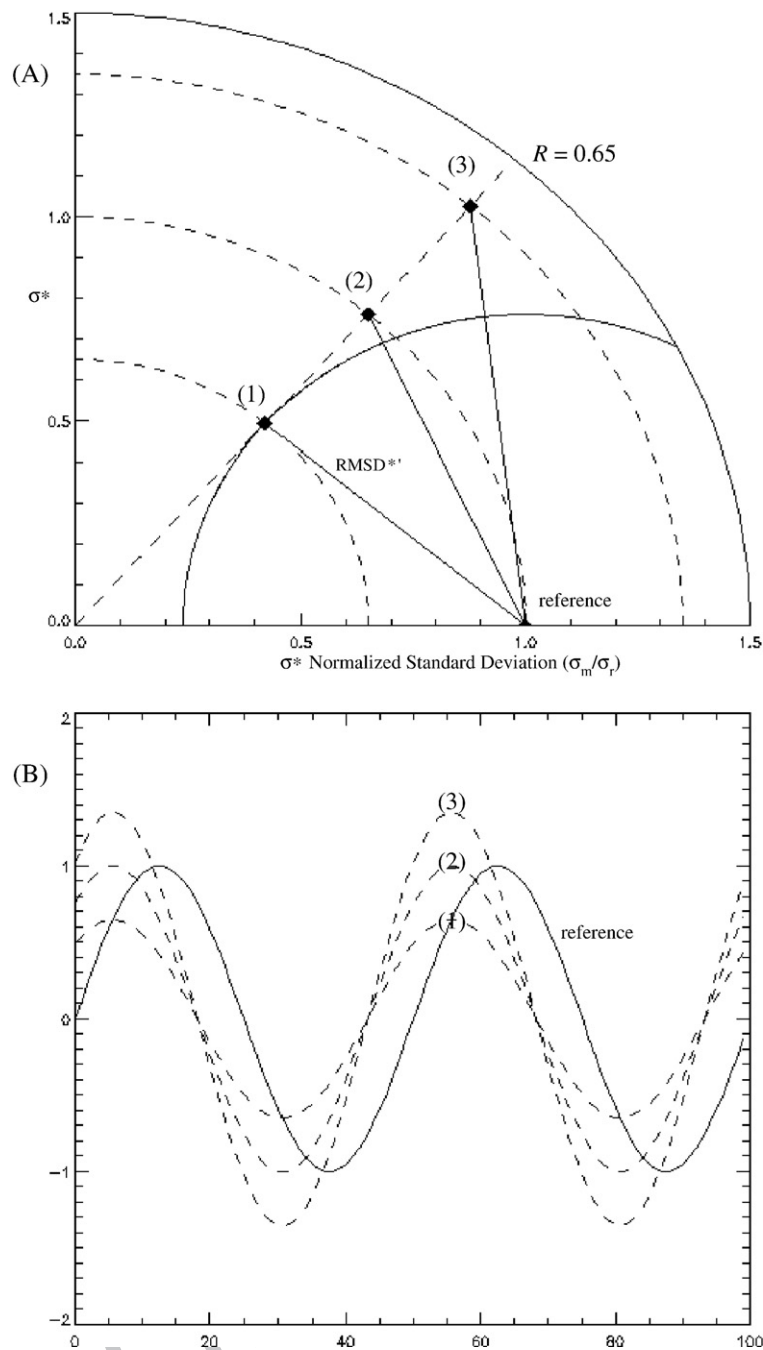


Fig. 9. (A) A Taylor diagram is shown for three model to reference field comparisons where $R=0.65$ and (1) $\sigma^*=0.65$, (2) $\sigma^*=1.0$, and (3) $\sigma^*=1.35$. An example of three sinusoidal waveforms and a reference field corresponding to the statistics in (A) is shown in panel (B).

660 values <1.0 are preferable. Clearly, if the two signals are out of
 661 phase, then reduction in the model variance to a threshold value
 662 diminishes the total RMSD value. However, if the goal of the
 663 investigation is to independently move R and σ^* as close to an
 664 ideal value of 1.0 as is possible then it may be inappropriate to
 665 use the total or unbiased RMSD as a model validation metric.

666 This is an important point since many model and observa-
 667 tion comparison exercises may involve RMSD-based
 668 metrics. For example, Wallhead et al. (submitted for publica-

tion) use the term “skillful” to refer to model predictions that
 669 minimize mean-square differences. Sheng and Kim (sub-
 670 mitted for publication) use RMSD metrics and Taylor dia-
 671 grams as part of their water quality model evaluation scheme.
 672 Smith et al. (submitted for publication) use an RMSD-based
 673 cost function as part of a data assimilation scheme. Indeed,
 674 RMSD-based metrics of model performance are likely to con-
 675 tinue to be used in a wide variety of contexts and investiga-
 676 tors should at least be cognizant of how RMSD-based functions or
 677

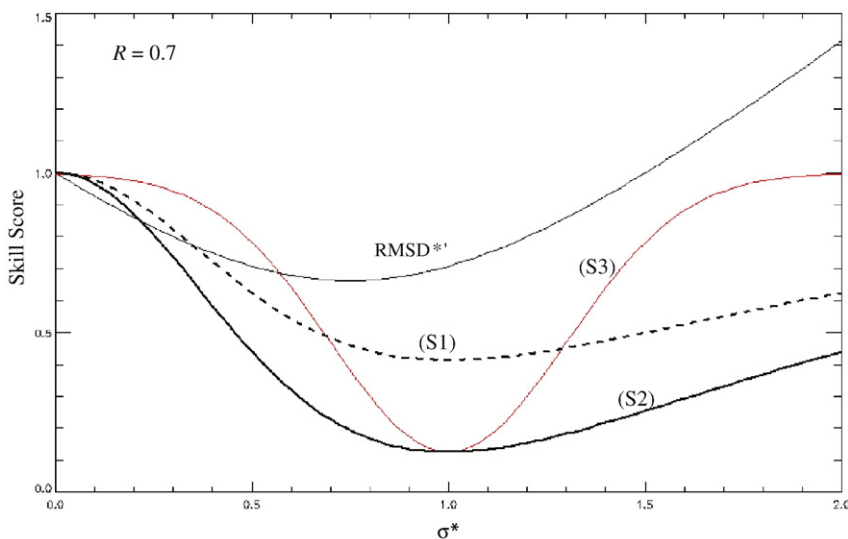


Fig. 10. The unbiased RMSD and skill scores S1–S3 are shown for $R=0.7$ and σ^* over the range $[0, 2]$.

678 skill scores quantify mismatches in variance when correlation
679 coefficients are less than unity.

680 Alternative metrics of model skill (skill scores) have been
681 proposed (Murphy and Epstein, 1989; Taylor, 2001), such as:

$$S1 = 1.0 - \frac{2(1+R)}{(\sigma^* + 1/\sigma^*)^2} \quad (16)$$

683 and

$$S2 = 1.0 - \frac{(1+R)^4}{(\sigma^* + 1/\sigma^*)^4} \quad (17)$$

685 The prevailing convention is to have the skill score range
686 between 0.0 (for poor skill) and 1.0 (for superior skill). This
687 convention is reversed here since our objective is to build a
688 summary skill target diagram similar to the one developed in
689 Section 3.2.

690 An important feature to consider is how these potential
691 skill scores proportionally penalize underestimates or over-
692 estimates of the standard deviation. For example, given a
693 constant R value of 0.7, the normalized and unbiased RMSD,
694 S1, and S2 are shown for $0.0 \leq \sigma^* \leq 2.0$ in Fig. 10. Minimum
695 skill scores occur where $\sigma^* = 1$, consistent with our stated skill
696 score convention. However, S1 and S2 appear to penalize
697 underestimates of the variance more than proportional over-
698 estimates, and are thus opposite of the $RMSD^*_{\lambda}$ statistic that
699 rewards variance underestimates. A potential alternative to
700 these measures is a Gaussian function that penalizes propor-
701 tional overestimates and underestimates of σ^* equally over
702 the range $[0, 2]$. Multiplication by a scaled correlation score
703 may then constitute a measure of model skill:

$$S3 = 1.0 - \left(e^{-\frac{(\sigma^* - 1.0)^2}{0.18}} \right) \left(\frac{1+R}{2} \right) \quad (18)$$

706 This measure of skill may now be incorporated into a
707 diagram similar to the one developed in the previous section.
708 Here, however, the emphasis is on the comparison of one

709 model to another more than the misfit between the model
710 and the data. Accordingly, a relative measure of bias may be
711 given as:

$$B_m = \frac{B_i}{|\text{Max}\{B_{i=1,2,3 \dots n}\}|} \quad (19)$$

712 that is, the maximum normalized bias of the i th model exe-
713 cution is its bias divided by the maximum magnitude bias
714 from the total set of n model to data comparisons.

715 If B_m serves as the Y-axis and S3 times the sign of the
716 standard deviation difference (σ_d) serves as the X-axis, then
717 the resulting skill target diagram renders distances from the
718 origin that are proportional to:
719

$$ST = \sqrt{B_m^2 + S3^2} \quad (20)$$

720 The contrast between the ST score and the total RMSD is that
721 the skill score does not reward underestimates of the variance
722 for correlation values less than one. Markers for the skill target
723 diagram are based on the percentile ST score of the models. For
724 example, in this case the mean ST score (\overline{ST}) is 0.51 and the
725 standard deviation (σ_{ST}) is 0.28, thus the 90th percentile
726 (assuming a normal score probability density function and
727 recalling our skill convention rewards low scores instead of
728 high scores) corresponds to $\overline{ST} - 1.28 \sigma_{ST}$ or $ST = 0.15$. A similar
729 marker for the 50th percentile ($ST = \overline{ST}$) is shown on Fig. 11. In
730 this case, the most skillful simulation (point 2, Fig. 11) is yet
731 again different from the minimum total RMSD simulation
732 (point 1, Fig. 11).
733

734 The discrepancy between minimum skill and RMSD scores is
735 exaggerated for the phytoplankton absorption field (Fig. 12).
736 The minimum unbiased RMSD score, as would appear to be the
737 best fit in a Taylor diagram, is also indicated (point 3, Fig. 12).
738 These three model fields are presented against the reference
739 field in Fig. 13. Evidently, the minimum unbiased RMSD model
740 field is unacceptable due to the large positive bias. In contrast,
741 the minimum RMSD (point 1, Fig. 12) and superior skill model
742 fields (point 2; Fig. 12) are less biased but are out of phase with
743

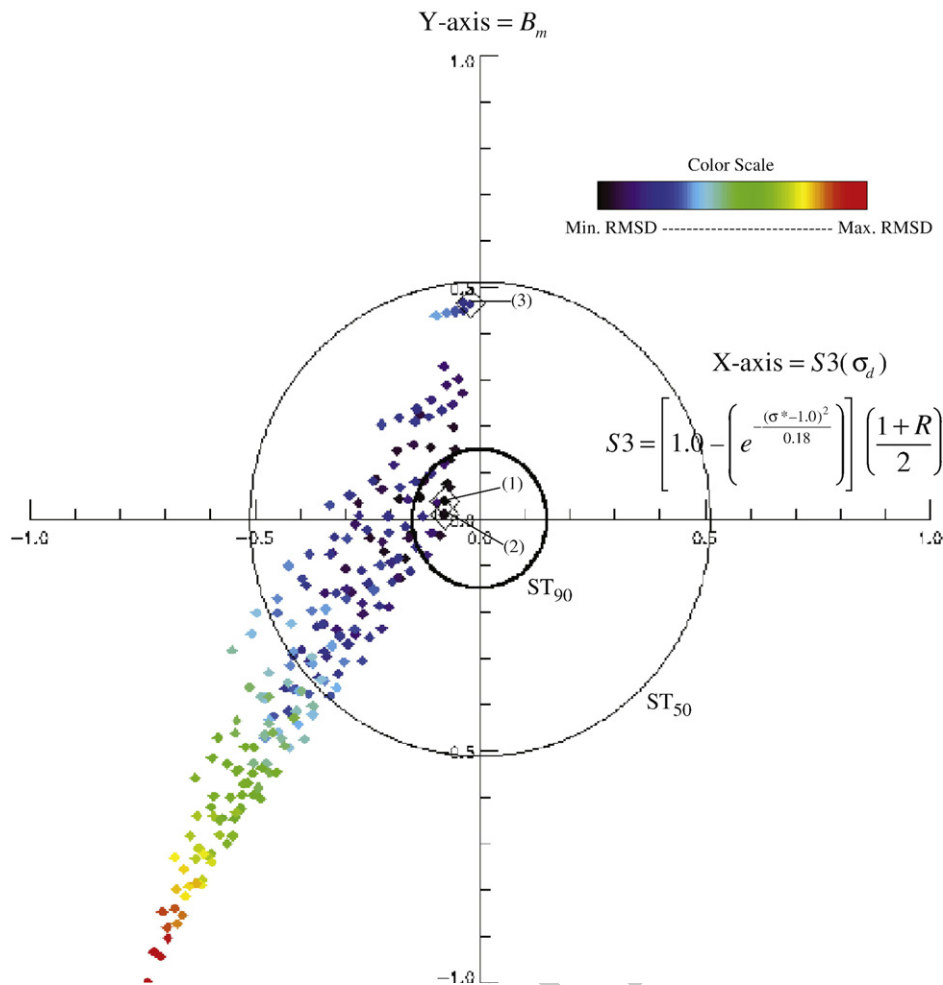


Fig. 11. Skill target diagram for model to reference chlorophyll-*a* field comparisons. The minimum total RMSD (1), minimum skill score (2), and minimum unbiased RMSD (3) are indicated on the plot. The markers indicate the 50th and 90th percentile total skill scores (ST) for the total set of model to reference comparisons, as explained in the text. The X-axis is the S3 skill score multiplied by the sign of the standard deviation difference. The Y-axis is the maximum normalized bias. The color scale indicates the total RMSD values.

743 the reference field by several months (Fig. 13). All three results
744 provide information potentially useful to the investigator; other
745 parameters may potentially be adjusted to either reduce the
746 phase error for fields (1) and (2), or the bias may be reduced in
747 (3), which is better correlated with the reference field. The
748 salient point to be made here, however, is that for multiple
749 model executions the skill target diagram may identify poten-
750 tial contrasts between minimum RMSD and other measures of
751 model skill.

752 4. Discussion

753 An important point mentioned elsewhere in this special
754 volume (Stow et al., submitted for publication) is worthy of
755 reiteration here: different statistical quantities (i.e., skill
756 metrics) may capture different aspects of model performance,
757 and a thorough assessment of model skill may require use of
758 multiple types of skill metrics simultaneously. Accordingly, it
759 is important to recognize the relationships that exist between

various statistical quantities and how they represent related
but differentiable aspects of model performance. Linear cor-
relation coefficients and variance comparisons help to iden-
tify similarities of pattern, and they may be combined in a way
that is equivalent to the unbiased RMSD score (Eq. (7)), which
succinctly quantifies pattern agreement. In our example of a
one-dimensional time series, we related these aspects of
model performance to the similarity of phase and amplitude
between two time-dependent and sinusoidal-like patterns,
but this concept may be generalized to describe the shape
(such as the pattern of potential contour lines) of multidimen-
sional property fields.

Pattern agreement is an important aspect of model per-
formance, and there may be instances where this aspect is of
particular or exclusive concern to the investigator. For exam-
ple, Li et al. (2007) use Taylor diagrams to compare modeled
and observed distributions of soil moisture and precipitation.
Since the average values from the simulations were adjusted
to agree with observed averages, the pattern information was

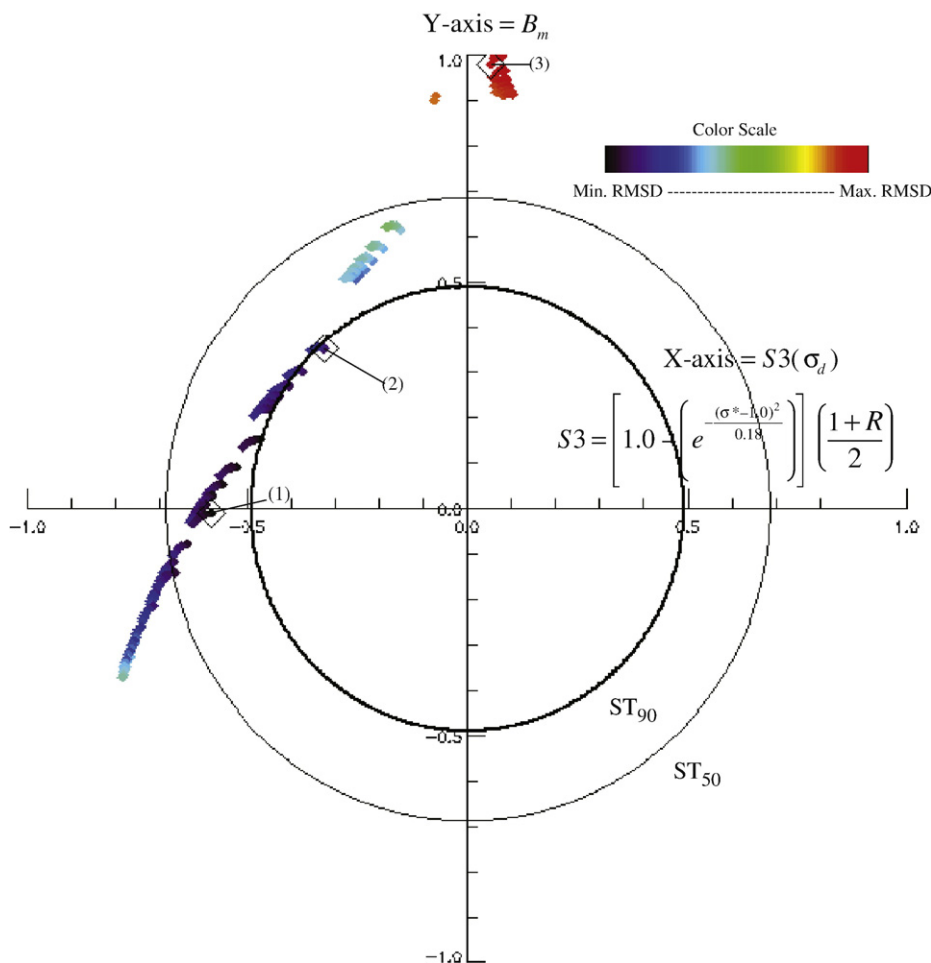


Fig. 12. Skill target diagram for model to reference phytoplankton absorption field comparisons. The minimum total RMSD (1), minimum skill score (2), and minimum unbiased RMSD (3) are indicated on the plot. The markers indicate the 50th and 90th percentile total skill scores (ST) for the total set of model to reference comparisons, as explained in the text. The X-axis is the S3 skill score multiplied by the sign of the standard deviation difference. The Y-axis is the maximum normalized bias. The color scale indicates the total RMSD values.

779 the primary aspect of interest from their climate model's
 780 performance. In such cases, Taylor diagrams are useful skill
 781 assessment tools insofar as they provide summary informa-
 782 tion about how the linear correlation coefficient and the var-
 783 iance comparisons each contribute to the unbiased RMSD on
 a two-dimensional diagram. Indeed, the pattern information
 may often be the primary area of interest for many climate
 model studies.

784 Nevertheless, in cases where the magnitude of the model
 785 results are not adjusted *a posteriori*, the usefulness of the Taylor
 786 diagram (and the statistical quantities it summarizes) as a skill
 787 assessment tool may be incomplete since it often provides no
 788 information about other aspects of model performance such as
 789 the bias (the comparison of mean values) or the total RMSD (a
 790 metric for overall model and data agreement). One way to
 791 remedy this omission is to modify Taylor diagrams via the ad-
 792 dition of a color dimension indicating the magnitude of either
 793 the bias or the total RMSD. An example of this style of modi-
 794 fication is given here and has been previously shown elsewhere
 795 (Orr, 2002).

796 More generally, however, information about the bias intro-
 797 duces the aspect of scale or magnitude to the model skill
 798 assessment process. For example, two surface chlorophyll
 799 fields may have a perfect correlation score and identical
 800 variances but the model field may still be an order of
 801 magnitude larger than the observations. This would suggest
 802 that too much nitrogen or carbon, for example, resides within
 803 the phytoplankton compartment and the ecosystem model
 804 may be inappropriately parameterized or structurally inade-
 805 quate. In many ocean ecosystem (or biogeochemical) model
 806 applications, the time-dependent flux of materials from one
 807 reservoir to another may be constrained by the magnitude of
 808 the observations, rather than merely the pattern information.
 809 This is particularly pertinent to the biological aspects of
 810 coupled models because the overall magnitude of biological
 811 productivity is a critical aspect of ecosystem function. Fur-
 812 thermore, while the unbiased RMSD may effectively quantify
 813 pattern agreement, it is seldom used as a metric for overall
 814 model and data agreement, whereas the total RMSD is more
 815 frequently applied to this task.

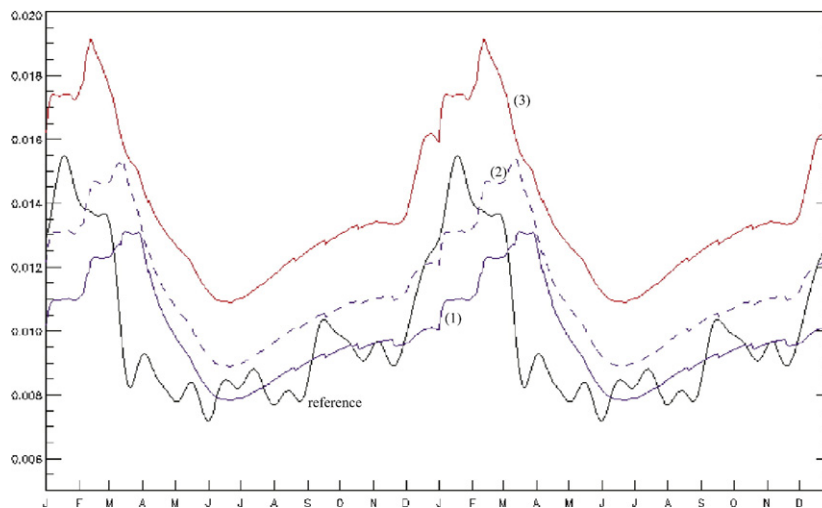


Fig. 13. The model and reference fields are plotted for the results indicated in Fig. 12: the minimum total RMSD (1), minimum skill score (2), and minimum unbiased RMSD (3; red). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

819 For these reasons, we have developed the target diagram, a
 820 Cartesian coordinate plot that provides summary information
 821 about how the magnitude and sign of the bias and the pattern
 822 agreement (unbiased RMSD) each contribute to the total
 823 RMSD magnitude. Markers may be added to the diagram in
 824 order to: (1) help identify limits based upon the correlation
 825 coefficient; (2) provide an assessment of model performance
 826 compared to an observational average (marker M_0); and (3)
 827 indicate potential limits to model performance improvement
 828 when the average observational uncertainty has been esti-
 829 mated. The observational uncertainty marker creates a “bull’s-
 830 eye” for the target diagram that may very effectively com-
 831 municate the estimated limits of model performance to other
 832 investigators.

833 For example, in our sensitivity analysis of grazing para-
 834 meter selection, 216 model fields may be compared to three
 835 reference field categories for a total of 648 sets of model to
 836 reference field statistics. These may all be summarized on a
 837 single target diagram (Fig. 14). cursory inspection of this
 838 summary diagram reveals that phytoplankton absorption is
 839 the most sensitive field and CDM absorption is the least. The
 840 phytoplankton absorption field is also prone to a large posi-
 841 tive bias. The chlorophyll field appears to achieve the mini-
 842 mum magnitude for total difference statistics, but further
 843 improvement would be within the estimated range of average
 844 observational uncertainty.

845 To be sure, the purpose of both the Taylor and target dia-
 846 grams is to compactly summarize statistical quantities that
 847 serve to aid in the skill assessment of model performance. The
 848 utility of either approach is dependent upon the aspects of
 849 model performance the metrics they summarize adequately
 850 capture. For the specific application to ocean ecosystem model-
 851 ing, we suggest that target diagrams may better summarize the
 852 overall agreement between model and data since aspects of
 853 pattern agreement and magnitude (bias) are given equal weight
 854 and one may clearly visualize how they each contribute to the
 855 total RMSD.

856 It would be inappropriate, however, to suggest that skill
 857 assessment must always be implicitly synonymous with finding

858 the lowest RMSD value amongst an ensemble of model results
 859 or an acceptably low RMSD values for a single model result. A
 860 potential deficiency in both the Taylor and target diagrams
 861 stems directly from a peculiarity of the RMSD metrics: the
 862 RMSD values may improve for correlations less than unity
 863 ($R < 1.0$) where the normalized standard deviation is equal to the
 864 correlation ($\sigma^* = R$) instead of an ideal value of one ($\sigma^* = 1.0$).

865 Another way to conceive of this behavior: if the correlation
 866 between a modeled and observed field is imperfect, i.e., in some
 867 areas the modeled values increase where or when the observed
 868 values decrease, then the average magnitude of this misfit may
 869 be reduced by diminishing the observed field’s variance (as-
 870 suming the bias is not a significant source of mismatch). For
 871 example, suppose a three-dimensional coupled model of phy-
 872 toplankton growth and ocean circulation appears to adequately
 873 reproduce the observed details of chlorophyll patterns within a
 874 mesoscale eddy, only the eddy is in the wrong location when
 875 compared to the observations (a common type of mismatch for
 876 coupled models since modeled velocity fields are imperfect and
 877 advection is a time-integrative process). Given this spatial
 878 mismatch, the RMSD-based metrics of model/data misfit may
 879 improve if the details (i.e., the variance) of the modeled
 880 chlorophyll field are diminished or smoothed over. Would the
 881 investigator prefer a blurred modeled field over the one where
 882 the exclusive source of model/data disagreement appears to be
 883 dislocation?

884 This circumstance may be clearly demonstrated using
 885 satellite ocean color patterns from areas of complex mesoscale
 886 variability, such as Moderate Resolution Imaging Spectro-
 887 radiometer data for the Mozambique Channel off the south-
 888 west coast of Madagascar (Fig. 15A). The complex pattern of
 889 apparent surface chlorophyll within mesoscale eddies and
 890 fronts (Fig. 15A) may potentially be mimicked by a coupled
 891 model, but imperfectly so with respect to spatiotemporal
 892 agreement. We approximate this kind of disagreement by
 893 reversing the array order (Fig. 15B) such that the hypothetical
 894 modeled field is effectively a mirror image of the data. The
 895 means and variances of the two fields are identical, but the
 896 correlation between them is quite low ($R = 0.09$) and this

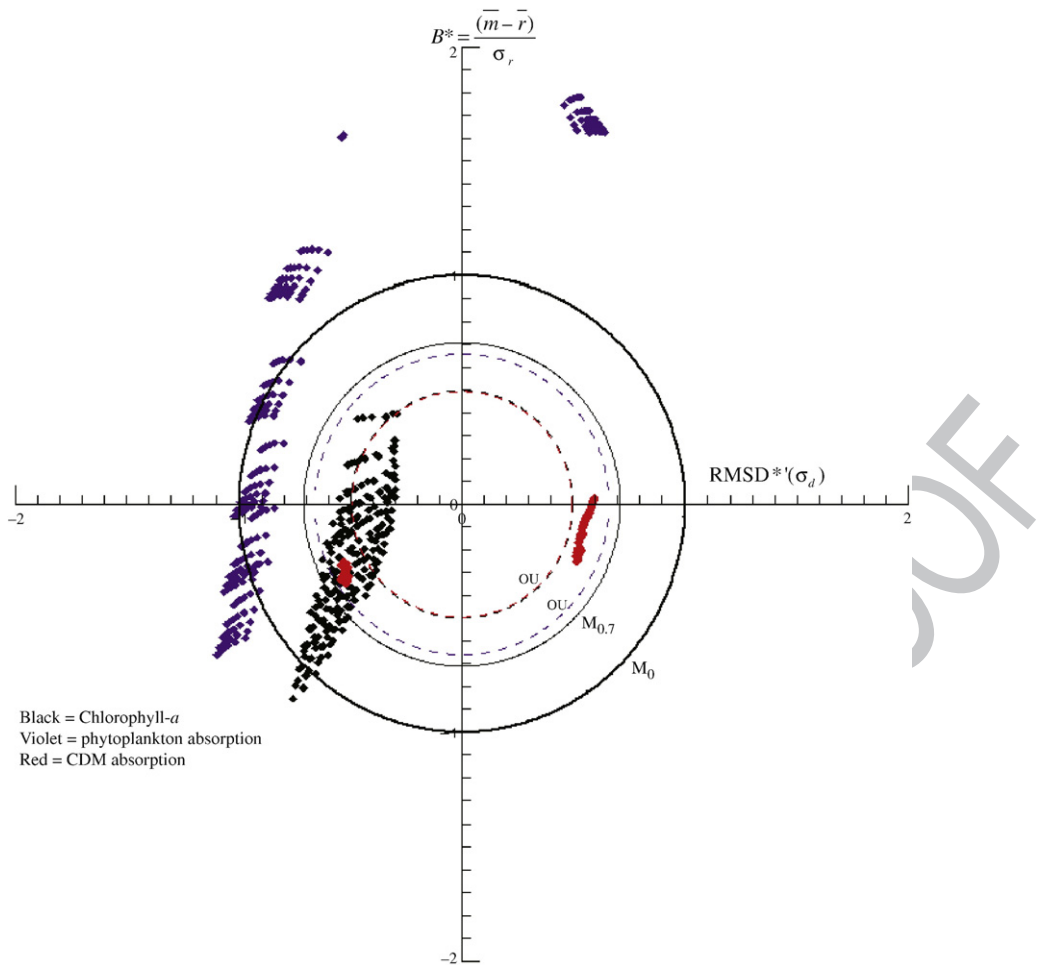


Fig. 14. Summary target diagram for all three types of model to reference field comparisons: chlorophyll- a (black), phytoplankton absorption (violet), and CDM absorption (red). The dashed lines indicate the estimated observational uncertainty (OU) threshold (corresponding to the field color). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

897 results in high RMSD scores ($\text{RMSD}^* = \text{RMSD} = 1.35$). These
 898 scores may be artificially improved by simply reducing the
 899 variance of the hypothetical model field (Fig. 15C) until the
 900 threshold criterion $\sigma^* = R$ is met. As a result of this procedure,
 901 complex spatial details of the modeled chlorophyll field have
 902 been significantly diminished (Fig. 15B and C) yet the RMSD
 903 scores have certainly improved ($\text{RMSD}^* = 0.99$). Another way
 904 to demonstrate this property of RMSD-based metrics is to
 905 begin with the original field (Fig. 15A) and simply apply a large
 906 smoothing filter (Fig. 15D). Of the three hypothetical modeled
 907 fields (Fig. 15B,C, and D), one may be inclined to select B as the
 908 most skillful, though RMSD scores run contrary to this
 909 inclination.

910 Thus there are indeed cases where a distinction may be
 911 appropriately made between reducing RMSD statistics and
 912 increasing model skill. An alternative skill scoring system and
 913 skill target diagram was developed and presented for such a
 914 contingency. The advantage of this system is that for $R < 1.0$
 915 the minimum value skill score instead occurs where $\sigma^* = 1.0$.
 916 In our example, the S3 skill score, Eq. (18), would indicate that
 917 field (B) is indeed the most skillful (Fig. 15). There are

potentially many other creative ways to combine correlations,
 918 variances, and other metrics into composite skill scores that
 919 have properties distinctly different from RMSD-based met-
 920 rics. Our intent is not to promote a specific solution but,
 921 rather, to point out that a contradiction may arise between
 922 minimum RMSD scores and other potential definitions of
 923 model skill.

924 In summary, model skill assessment ultimately requires
 925 specification about which quantitative metrics should be
 926 applied and how they should be interpreted to constitute
 927 “good” or “bad” model performance. The “skill” portion of
 928 skill assessment may be mathematically defined, but the
 929 “assessment” will invariably rely upon the value judgments of
 930 the investigator. Our analysis has focused upon some widely
 931 known statistical quantities (linear correlation coefficients,
 932 means, and variances) and ways that they may be combined
 933 mathematically and graphically to describe RMSD-based
 934 measures of model/data misfit. Taylor diagrams are polar
 935 coordinate plots that focus upon pattern agreement, whereas
 936 the target diagrams developed here summarize both the
 937 aspects of pattern agreement and magnitude (bias) and how
 938

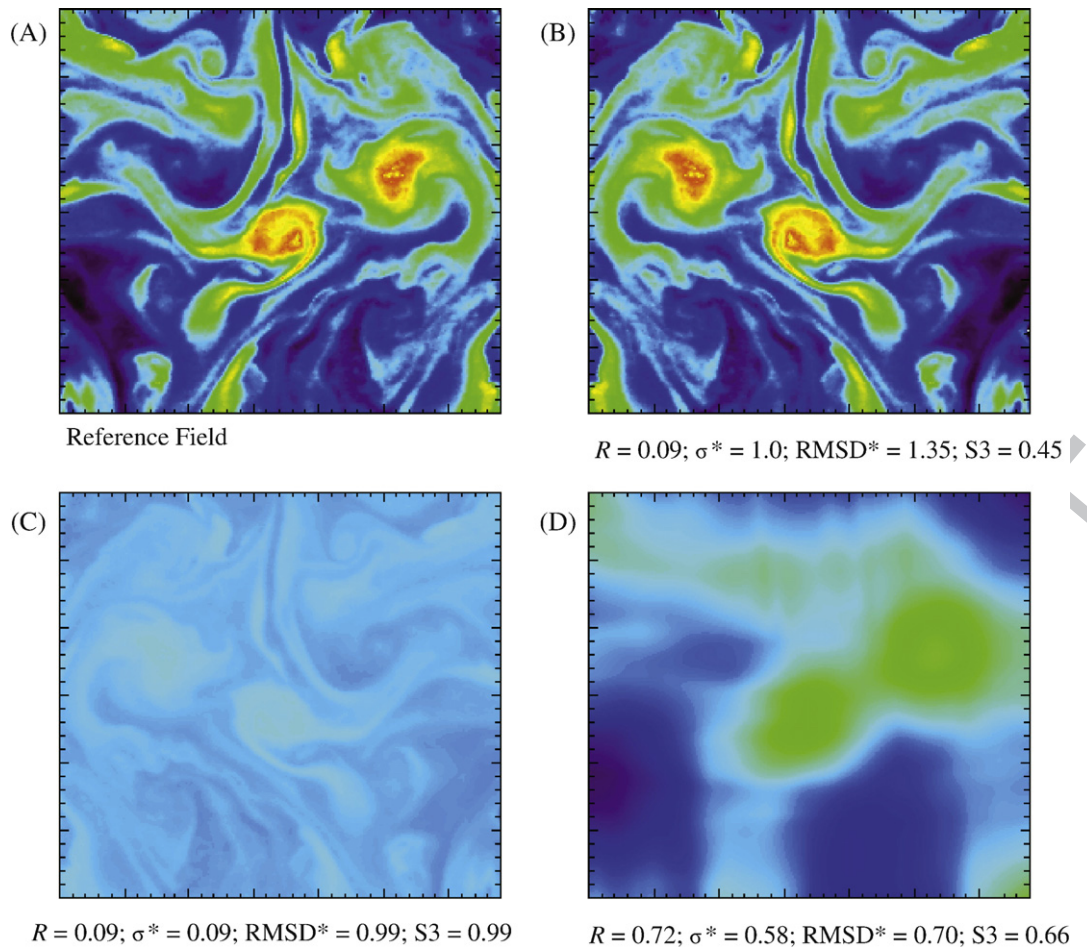


Fig. 15. A pattern of ocean color data is shown in panel A (surface chlorophyll fields; Moderate Resolution Imaging Spectroradiometer image 25 July 2007; data provided by NASA from their website at <http://oceancolor.gsfc.nasa.gov/>). To make a hypothetical model field wherein the misfit arises exclusively from spatial incoherence, the data array in (A) was reversed and is shown in panel (B) as a hypothetical modeled field. The resulting correlation is low but the mean and variance are the same. The field in panel (B) was further manipulated so that the normalized standard deviation (σ^*) is equal to the correlation coefficient ($\sigma^* = R$). This field is shown in panel (C). As a final comparison, the field in panel (A) was smoothed using a moving average filter. The correlation (R), normalized standard deviation (σ^*), normalized total root-mean-square difference (RMSD^*), and skill score ($S3$) are shown beneath each panel for the comparison to the reference field (A). Panel (D) has the lowest RMSD^* score and panel (B) has the lowest skill score.

939 they each contribute to the total RMSD, a common metric of
 940 overall model/data agreement. Investigators should be cogni-
 941 zant of the aspects of model performance summarized by
 942 each of these aforementioned statistical and graphical ap-
 943 proaches before making claims of “model validation.” Further-
 944 more, both methods presume that RMSD-based metrics are
 945 sufficient criteria upon which to base model skill assessments,
 946 and this may not always be the case.

947 Acknowledgements

948 This research is a contribution to the Naval Research Labo-
 949 ratory 6.1 project, “Coupled Bio-Optical and Physical Processes
 950 in the Coastal Zone” under program element 61153N sponsored
 951 by the Office of Naval Research. This research was supported by
 952 the National Research Council’s Post-Doctoral Research Asso-
 953 ciateship Program, and partially supported by the Office of
 954 Naval Research, grant number N0001405WX20735. Paul
 955 Martinolich provided assistance with SeaWiFS data processing

and C. N. Barron and Clark Rowley provided assistance with the 956
 MODAS system. We also would like to thank two anonymous 957
 reviewers whose helpful comments certainly improved this 958
 manuscript. 959

References 960

- Allen, J.I., Somersfield, P.J., Gilbert, F.J., 2007. Quantifying uncertainty in high- 961
 resolution coupled hydrodynamic-ecosystem models. *Journal of Marine* 962
Systems 64, 3–14. 963
 Bailey, S.W., Werdell, P.J., 2006. A multi-sensor approach for the on-orbit 964
 validation of ocean color satellite data products. *Remote Sensing of* 965
Environment 102, 12–23. 966
 Bretherton, F.P., Davis, R.E., Fandry, C.B., 1976. A technique for objective 967
 analysis and design of oceanographic experiments applied to MODE-73. 968
Deep-Sea Research 23, 559–582. 969
 Bricaud, A., Claustre, H., Ras, J., Oubelkheir, K., 2004. Natural variability of 970
 phytoplankton absorption in oceanic waters: influence of the size 971
 structure of algal populations. *Journal of Geophysical Research* 109, 972
 C11010. doi:10.1029/2004JC002419. 973
 Fox, D.N., Teague, W.J., Barron, C.N., Carnes, M.R., Lee, C.M., 2002. The Modular 974
 Ocean Data Assimilation System (MODAS). *Journal of Atmospheric and* 975
Oceanic Technology 19, 240–252. 976

- 977 Franks, P.J.S., Chen, C., 2001. A 3-D prognostic numerical model study of the
978 Georges bank ecosystems. Part II: biological-physical model. *Deep-Sea*
979 *Research II* 48, 457–482.
- 980 Friedrichs, M.A.M., Dusenberry, J., Anderson, L.A., Armstrong, R.A., Chai, F.,
981 Christian, J.R., Doney, S.C., Dunne, J., Fujii, M., Hood, R., McGillicuddy, D.J.,
982 Moore, J.K., Schartou, M., Spitz, Y.H., Wiggert, J.D., 2007. Assessment of
983 skill and portability in regional marine biogeochemical models: role of
984 multiple planktonic groups. *Journal of Geophysical Research* 112, C08001.
985 doi:10.1029/2006JC003852.
- Q1 986 Friedrichs, M.A.M., Carr, M.-E., Scardi, M., Barber, R., submitted for publica-
987 tion. Assessing the uncertainties of model estimates of primary
988 productivity in the tropical Pacific Ocean. *Journal of Marine Systems*.
- 989 Gregg, W.W., Ginoux, P., Schopf, P.S., Casey, N.W., 2003. Phytoplankton and
990 iron: validation of a global three-dimensional ocean biogeochemical
991 model. *Deep-Sea Research II* 50, 3143–3169.
- 992 Gruber, N., Frenzel, H., Doney, S.C., Marchesiello, P., McWilliams, J.C., Moisan,
993 J.R., Oram, J.J., Plattner, G.-K., Stolzenbach, K.D., 2006. Eddy-resolving
994 simulation of plankton ecosystem dynamics in the California Current
995 System. *Deep-Sea Research I* 53, 1483–1516.
- 996 Holt, J.T., Allen, J.L., Proctor, R., Gilbert, F.G., 2005. Error quantification of a
997 high resolution coupled hydrodynamic-ecosystem coastal ocean
998 model: Part 1. Model overview and hydrodynamics. *Journal of Marine*
999 *Systems* 57, 167–188.
- 1000 Ivlev, V.S., 1961. *Experimental Ecology of the Feeding of Fishes*. Yale
1001 University Press, New Haven, Connecticut. 302 pp.
- 1002 Jochens, A.E., DiMarco, S.F., Nowlin Jr., W.D., Reid, R.O., Kennicutt II, M.C.,
1003 2002. Northeastern Gulf of Mexico Chemical Oceanography and Hydro-
1004 graphy Study: Synthesis Report. Technical Report, U.S. Department of the
1005 Interior, Minerals Management Service, Gulf of Mexico OCS Region, New
1006 Orleans, Louisiana. 586 pp.
- 1007 Jolliff, J.K., Kindle, J.C., 2007. Naval Research Laboratory Ecological-Photo-
1008 chemical-Bio-Optical-Numerical Experiment (Neptune) Version 1: a
1009 portable, flexible modeling environment designed to resolve time-
1010 dependent feedbacks between upper ocean ecology, photochemistry,
1011 and optics. NRL Technical Memorandum, NRL/MR/7330-07-9026, Naval
1012 Research Laboratory, Stennis Space Center, Mississippi. 49 pp., [http://](http://stinet.dtic.mil/)
1013 stinet.dtic.mil/.
- 1014 Kindle, J.C., DeRada, S., Arnone, R.A., Shulman, I., Penta, B., Anderson, S., 2005.
1015 Near real-time depiction of the California Current System. American
1016 Meteorological Society Sixth Conference on Coastal Atmospheric and
1017 Oceanic Prediction and Processes, San Diego, CA.
- 1018 Lee, Z.P., Carder, K.L., Arnone, R.A., 2002. Deriving inherent optical properties
1019 from water color: a multiband quasi-analytical algorithm for optically
1020 deep waters. *Applied Optics* 41, 5755–5772.
- 1021 Li, H., Robok, A., Wild, M., 2007. Evaluation of Intergovernmental Panel on
1022 Climate Change Fourth Assessment soil moisture simulations for the
1023 second half of the twentieth century. *Journal of Geophysical Research*
1024 112, D06106. doi:10.1029/2006JD007455.
- 1025 McClain, C., Hooker, S., Feldman, G., Bontempi, P., 2006. Satellite data for
1026 ocean biology, biogeochemistry, and climate research. *EOS Transactions,*
1027 *American Geophysical Union* 87, 337.
- 1028 Millan-Nunez, E., Sieracki, M.E., Millan-Nunez, R., Lara-Lara, J.R., Gaxiola-
1029 Castro, G., Trees, C.C., 2004. Specific absorption coefficient and phyto-
1030 plankton biomass in the southern region of the California Current. *Deep-*
1031 *Sea Research II* 51, 817–826.
- 1032 Murphy, A.H., Epstein, E.S., 1989. Skill scores and correlation coefficients in
1033 model verification. *Monthly Weather Review* 117, 572–581.
- 1034 Nash, J.E., Sutcliffe, J.V., 1970. River flow forecasting through conceptual
1035 models, Part 1 – A discussion of principles. *Journal of Hydrology* 10,
1036 282–290.
- 1037 O'Reilly, J.E., Maritorea, S., Mitchell, B.G., Siegal, D.A., Carder, K.L., Garver,
1038 S.A., Kahru, M., McClain, C., 1998. Ocean color algorithms for SeaWiFS.
1039 *Journal of Geophysical Research* 103, 24937–24953.
- 1040 Orr, J.C., 2002. Global Ocean Storage of Anthropogenic Carbon (GOSAC). Final
1041 Report (December 1, 1997 to March 31, 2001). EC Environmental and
1042 Climate Programme (Contract ENV4-CT97-0495). IPSL/CNRS, France.
1043 116 pp.
- 1044 Pacanowski, R.C., Philander, S.G.H., 1981. Parameterization of vertical mixing
1045 in numerical models of the tropical oceans. *Journal of Physical Oceano-*
1046 *graphy* 11, 1443–1451.
- 1047 Raick, C., Alvera-Azcarate, A., Barth, A., Brankart, J.M., Soetaert, K., Gregoire,
1048 M., 2007. Application of a SEEK filter to a 1D biogeochemical model of the
1049 Ligurian Sea: twin experiments and real in-situ data assimilation.
1050 *Journal of Marine Systems* 65, 561–583.
- 1051 Sheng, P., Kim, T., submitted for publication. Skill assessment of an integrated
1052 modeling system for shallow coastal and estuarine ecosystems. *Journal*
1053 *of Marine Systems*.
- 1054 Smith, K.W., McGillicuddy, D.J. Jr., Lynch, D.R., submitted for publication.
1055 Parameter estimation using an ensemble smoother: the effect of the
1056 circulation in biological estimation. *Journal of Marine Systems*.
- 1057 Stow, C.A., Jolliff, J.K., McGillicuddy, D.J. Jr., Doney, S.C., Allen, J.L., Rose, K.A.,
1058 Wallhead, P., submitted for publication. Skill assessment for coupled
1059 biological/physical models of marine systems. *Journal of Marine Systems*.
- 1060 Stow, C.A., Roessler, C., Borsuk, M.E., Bowen, J.D., Reckhow, K.H., 2003. A
1061 comparison of estuarine water quality models for TMDL development in
1062 the Neuse River Estuary. *Journal of Water Resources Planning and*
1063 *Management* 129, 307–314.
- 1064 Taylor, K.E., 2001. Summarizing multiple aspects of model performance in a
1065 single diagram. *Journal of Geophysical Research* 106, 7183–7192.
- 1066 Wallhead, P.J., Martin, A.P., Srokosz, M.A., Franks, P.J.S., submitted for
1067 publication. Predicting the bulk plankton dynamics of Georges Bank:
1068 model skill assessment. *Journal of Marine Systems*.
- 1069 Walsh, J.J., Weisberg, R.H., Dieterle, D.A., He, R., Darrow, B.P., Jolliff, J.K., Lester,
1070 K.M., Vargo, G.A., Kirkpatrick, G.J., Fanning, K.A., Sutton, T.T., Jochens, A.E.,
1071 Biggs, D.C., Nababan, B., Hu, C., Muller-Karger, F.E., 2003. The phyto-
1072 plankton response to intrusions of slope water on the West Florida shelf:
1073 models and observations. *Journal of Geophysical Research* 108 (C6), 3190.
1074 doi:10.1029/2002JC001406.