# Challenges associated with modeling low oxygen waters in Chesapeake Bay: a

# **3 multiple model comparison**

4 I. D. Irby<sup>1</sup>, M. A. M. Friedrichs<sup>1</sup>, C. T. Friedrichs<sup>1</sup>, A. J. Bever<sup>2</sup>, R. R. Hood<sup>3</sup>, L. W.

5 J. Lanerolle<sup>4,5</sup>, M. Li<sup>6</sup>, L. Linker<sup>7</sup>, M. E. Scully<sup>8</sup>, K. Sellner<sup>9</sup>, J. Shen<sup>1</sup>, J. Testa<sup>6</sup>, H.

- 6  $Wang^3$ , P.  $Wang^{10}$ , and M. Xia<sup>11</sup>
- <sup>1</sup>Virginia Institute of Marine Science, College of William & Mary, P.O. Box 1346, Gloucester
   Point, VA 23062, USA
- <sup>9</sup> <sup>2</sup>Anchor QEA, LLC, 130 Battery Street, Suite 400, San Francisco, CA 94111, USA

<sup>3</sup>Horn Point Laboratory, University of Maryland Center for Environmental Science, P.O. Box
 775, Cambridge, MD 21613, USA

- <sup>4</sup>NOAA/NOS/OCS Coast Survey Development Laboratory, 1315 East–West Highway, Silver
   Spring, MD 20910, USA
- <sup>5</sup>ERT Inc., 14401 Sweitzer Lane Suite 300, Laurel, MD 20707, USA
- <sup>6</sup>Chesapeake Biological Laboratory, University of Maryland Center for Environmental Science,
   P.O. Box 38, Solomons, MD 20688, USA
- <sup>7</sup>US Environmental Protection Agency Chesapeake Bay Program Office, 410 Severn Avenue,
   Annapolis, MD 21403, USA
- <sup>8</sup>Woods Hole Oceanographic Institution, Applied Ocean Physics and Engineering Department,
   Woods Hole, MA 02543, USA
- <sup>9</sup>Chesapeake Research Consortium, 645 Contees Wharf Road, Edgewater, MD 21037, USA
- 22 <sup>10</sup>VIMS/Chesapeake Bay Program Office, 410 Severn Avenue, Annapolis, MD 21403, USA
- 23 <sup>11</sup>Department of Natural Sciences, University of Maryland Eastern Shore, MD, USA

24

- 25
- 26

#### 28 Abstract

29 As three-dimensional (3-D) aquatic ecosystem models are used more frequently for 30 operational water quality forecasts and ecological management decisions, it is important 31 to understand the relative strengths and limitations of existing 3-D models of varying 32 spatial resolution and biogeochemical complexity. To this end, two-year simulations of the Chesapeake Bay from eight hydrodynamic-oxygen models have been statistically 33 34 compared to each other and to historical monitoring data. Results show that although 35 models have difficulty resolving the variables typically thought to be the main drivers of 36 dissolved oxygen variability (stratification, nutrients, and chlorophyll), all eight models 37 have significant skill in reproducing the mean and seasonal variability of dissolved 38 oxygen. In addition, models with constant net respiration rates independent of nutrient 39 supply and temperature reproduced observed dissolved oxygen concentrations about as 40 well as much more complex, nutrient-dependent biogeochemical models. This finding 41 has significant ramifications for short-term hypoxia forecasts in the Chesapeake Bay, 42 which may be possible with very simple oxygen parameterizations, in contrast to the 43 more complex full biogeochemical models required for scenario-based forecasting. 44 However, models have difficulty simulating correct density and oxygen mixed layer 45 depths, which are important ecologically in terms of habitat compression. Observations 46 indicate a much stronger correlation between the depths of the top of the pycnocline and 47 oxycline than between their maximum vertical gradients, highlighting the importance of 48 the mixing depth in defining the region of aerobic habitat in the Chesapeake Bay when 49 low-oxygen bottom waters are present. Improvement in hypoxia simulations will thus 50 depend more on the ability of models to reproduce the correct mean and variability of the 51 depth of the physically driven surface mixed layer than the precise magnitude of the vertical density gradient. 52

- 53
- 54

55

56

57

58

59

60

61

62

#### 64 **1 Introduction**

65 Since the middle of the last century, anthropogenic impacts have dramatically decreased 66 water quality throughout the Chesapeake Bay (Boesch et al., 2001), one of the largest 67 estuaries in North America. Land-use change along with the industrialization and 68 urbanization of the Chesapeake Bay watershed have caused dramatic increases in nutrient 69 inputs to the Bay (Kemp et al., 2005), spurring additional primary production and 70 phytoplankton abundance (Harding and Perry, 1997). Because increased primary 71 production leads to more organic matter throughout the water column that is eventually 72 decomposed by bacteria, these increased nutrient inputs to the Bay have led to a 73 corresponding decrease in dissolved oxygen (DO) concentrations (Hagy et al., 2004). 74 Hypoxia, generally defined as the condition in which DO concentrations are below  $2mgL^{-1}$ , usually initiates seasonally in the northern portion of the Bay and expands 75 76 southward as summer develops (Kemp et al., 2009; Testa and Kemp, 2014). Although 77 hypoxia in the Chesapeake Bay has likely existed since European colonization (Cooper 78 and Brush, 1991, 1993), recent studies have highlighted an accelerated rise in the number and spatial extent of hypoxic, as well as anoxic (DO concentrations <0.2mgL<sup>-1</sup>), events 79 80 in the Bay since the 1950's, primarily attributed to increased anthropogenic nutrient input 81 (Hagy et al., 2004; Kemp et al., 2005; Gilbert et al., 2010). These impacts are likely to be 82 exacerbated by future climate change (Najjar et al., 2010; Meire et al., 2013; Harding et 83 al., 2015).

84 Interest in the ecological impacts of reduced DO concentrations has been elevated due to the observed proliferation of hypoxic events in the world's coastal oceans, creating vast 85 86 dead zone areas that compress suitable habitat for many marine species (Diaz, 2001; Diaz and Rosenberg, 2008; Pierson et al., 2009). Low-DO waters can greatly impact the 87 88 abundance and health of important ecological species, potentially resulting in suffocation 89 and major kills of fish, crabs, and shellfish (Breitburg, 2002; Ekau et al., 2010; Levin et al., 2009). While the presence of DO concentrations  $< 2 \text{ mg L}^{-1}$  have been shown to 90 91 decrease the abundance of fish larvae (Keister et al., 2000), some species can incur 92 negative health impacts and modify their behavior at significantly higher DO

concentrations (Vaquer-Sunyer and Duarte, 2008). DO concentrations of  $\sim 4 \text{ mg L}^{-1}$ 93 94 have been found to compress demersal fish habitat as fish seek out more oxygenated 95 waters (Buchheister et al., 2013). Zooplankton, a crucial food source for valuable species, 96 have also been found to exhibit changes in distribution and predation when subject to 97 large volumes of low-DO water, potentially leading to further impacts along the food 98 chain (Breitburg et al., 1997; Pierson et al., 2009). Invertebrates have similarly been 99 found to alter their behavior under low-DO conditions (Riedel et al., 2014). In the 100 Chesapeake Bay, multiple regulated fish species, such as striped bass and American shad, require oxygen restoration targets as high as  $5mgL^{-1}$  (USEPA, 2010). The greatest impact 101 of low DO concentrations spatially will depend on the specific living resource; however, 102 103 temporally, late spring to early fall is of most concern. As a result of the significant 104 ecological importance of oxygen on living resources in the Bay, DO concentrations are 105 used as a primary indicator in assessing water quality for Chesapeake Bay regulations 106 (Keisman and Shenk, 2013).

107 Improving the health of the Chesapeake Bay has become a priority for the Environmental 108 Protection Agency (EPA) along with the six states and Washington, DC that make up the 109 Bay watershed (Fig. 1), and together they have committed to utilizing a suite of 110 regulatory models to inform their management decisions (USEPA, 2010). The 111 Chesapeake Bay Program (CBP), a regional partnership that has led and directed the 112 restoration of the Chesapeake Bay since 1983, has undertaken an extensive modeling 113 effort of the Bay (Cerco and Cole, 1993; Cerco et al., 2002; Cerco and Noel, 2004, 2013). 114 This modeling system is being used by the CBP to estimate the aggregate effect of 115 changes in management practices, including land use, atmospheric deposition, animal 116 populations, and fertilizer and manure application. Recently, the modeling system has 117 been used to conduct scenario simulations to assess management actions needed to 118 achieve desired Bay water quality standards (USEPA, 2010). Ultimately this model was 119 used to establish a regulatory set of total maximum daily loads of nutrients and sediment 120 delivered from the watershed, with the goal of significantly improving water quality 121 throughout the Bay (USEPA, 2010).

122 Many 3-D hydrodynamic-oxygen models of varying complexity stemming from the 123 academic research community have also been used to simulate DO concentrations 124 throughout the Chesapeake Bay (Scully, 2010, 2013; Hong and Shen, 2013; Feng et al., 125 2015; Testa et al., 2014; Li et al., 2015). Bever et al. (2013) specifically demonstrated 126 that multiple models of varying complexity are able to generate skillful estimates of 127 hypoxic volume in the Bay. Some of these models are being used in the Bay to simulate 128 short-term and/or seasonal forecasts of DO conditions. Furthermore, some models are 129 also being used to generate scenario forecasts, or projections, that assess the impact of 130 changes in management practices on estuarine DO concentrations, in some cases taking 131 into account the impacts of future changes in climate.

132 As ecosystem and water quality models are increasingly used for operational forecasts as 133 well as scenario-based management decisions by the regulatory and academic research 134 communities, it is important to understand the relative strengths and limitations of 135 existing models of varying complexity. The ability to discern which variables must be 136 most accurately simulated in order to adequately reproduce the temporal and spatial 137 variability of Bay oxygen concentrations is a necessary prerequisite for fully 138 understanding how volumes of low-DO water are initiated and sustained within water 139 quality models. The utilization of multiple models can also inform projections by 140 providing independent confidence bounds for management decisions. To those ends, the 141 overarching goals of this research are to compare the relative skill of various three-142 dimensional (3-D) Chesapeake Bay models characterized by different levels of 143 biogeochemical complexity and spatial resolution, to better understand factors limiting 144 their ability to reproduce observed DO distributions, and to suggest approaches for the 145 continued improvement of these models.

#### 146 2 Methods

147 **2.1 Participating Chesapeake Bay models** 

148 Eight 3-D models were evaluated in this study (Table 1), each of which includes

- 149 hydrodynamic and DO components. Among the eight models, there are four different
- 150 hydrodynamic base models. Models B, C, D, F, and G utilize the Regional Ocean

151 Modeling System (ROMS; Shchepetkin and McWilliams, 2005; Haidvogel et al., 2008) 152 that employs a structured grid with sigma layers in the vertical dimension. Specifically, 153 Models B, C, and F use a ROMS implementation developed for the Chesapeake Bay 154 based on Xu et al. (2012; ChesROMS). Model D employs a ROMS implementation for 155 the Chesapeake Bay based on Li et al. (2005), while Model G uses the ROMS-based 156 Chesapeake Bay Operational Forecast System (CBOFS; Lanerolle et al., 2011). Models 157 A, E, and H each use a different hydrodynamic base model: the Curvilinear 158 Hydrodynamics in Three Dimensions model (CH3D; Cerco et al., 2010), the Finite-159 Volume Community Ocean Model (FVCOM; Jiang and Xia, 2015), and the 160 Hydrodynamic-Eutrophication Model – Fluid Dynamics Code (EFDC; Park et al., 1995; 161 Hong and Shen, 2012; Du and Shen, 2015), respectively. The only model that employs a 162 non-sigma vertical grid is Model A and the only model utilizing an unstructured 163 horizontal grid is Model E. While Model E contains 10 sigma vertical layers, all of the 164 other sigma grids use 20 layers. All of the grids vary in terms of their horizontal

resolution, with Models A and G utilizing the highest resolution horizontal grids.

166 These four hydrodynamic models are coupled to five different models used to simulate

167 DO (Table 1). Models A, B, C, D, and E utilize full biogeochemical models that include

as state variables various combinations of oxygen, phytoplankton, zooplankton, and

169 multiple inorganic and organic nutrients. Specifically, Models A and E employ a version

170 of the Integrated Compartment Model (ICM; Cerco et al., 2010; Jiang et al., 2015),

171 Model B uses the Estuarine Carbon Biogeochemistry model (ECB; Feng et al., 2015),

172 Model C uses the Biogeochemistry model (BGC; Brown et al., 2013), and Model D uses

the Row-Column AESOP model (RCA; Testa et al., 2014). In terms of food web

174 complexity the models vary considerably: Models B and C employ a single

175 phytoplankton group whereas Model D uses two phytoplankton groups, Model E uses

three, and Model A, the most complex of the participating models, uses five.

177 In contrast to the full biogeochemical models discussed above (Models A through E),

178 Models F, G, and H represent oxygen dynamics as simply as possible and therefore do

179 not utilize a full biogeochemical component. Rather, the models impose a biological

180 oxygen consumption rate that is model-specific, but constant in both space and time. This

- 181 component is referred to as a constant-respiration model (CRM). In this model, DO is
- 182 introduced to the estuary via the river and ocean boundaries and is set to saturation at the
- 183 estuarine surface. This constant-respiration oxygen parameterization (Scully, 2010) is
- simplistic, yet has been shown to adequately represent Chesapeake Bay oxygen dynamics
- 185 (Scully, 2010, 2013; Bever et al., 2013).
- 186 The major difference in forcing between the eight model implementations is that Models
- 187 A and B use riverine input derived from watershed models, whereas Models C–H used
- the measured flow from United States Geological Survey gauging stations, extrapolated
- using various techniques. Model A utilized the CBP's regulatory watershed model
- 190 (Shenk and Linker, 2013), while Model B utilized the Dynamic Land Ecosystem Model
- 191 (Yang et al., 2014, 2015; Tian et al., 2015). At the open boundary with the Atlantic
- 192 Ocean, Models B, C, D, F, G, and H utilize a sub-tidal elevation extrapolated from tidal
- stations on either side of the open boundary. Model E uses the TPXO tidal model, while
- 194 Model A uses a mix of observational and model forcing (Cerco et al., 2010). While
- 195 Model B utilizes wind forcing based on observations from the Thomas Point Light,
- 196 Models C through H use wind estimates from the North American Regional Reanalysis
- 197 (NARR).
- 198 The eight models used in this analysis have been developed for a variety of purposes.
- 199 Model A is a governmental regulatory model developed by the CBP that has been
- 200 extensively calibrated specifically to examine water quality issues in the Chesapeake Bay
- 201 (Cerco and Cole, 1993; Cerco and Noel, 2004, 2013; Cerco et al., 2010) and has been
- used in the development of the 2010 Chesapeake Bay Total Maximum Daily Load
- 203 (USEPA, 2010). The National Oceanic and Atmospheric Administration employs the
- 204 hydrodynamic component of Model F for operational forecasts of a variety of physical
- 205 estuarine parameters for the Chesapeake Bay
- 206 (http://www.tidesandcurrents.noaa.gov/ofs/cbofs/cbofs.html). The other six models are
- 207 academic models used in diverse research efforts focused on the Chesapeake Bay but not
- 208 necessarily specifically on DO dynamics.
- Finally, a ninth model is calculated as the mean of the results from the eight models

210 described above, and is referred to here as Model Mean, or Model M.

#### 211 **2.2** Available Chesapeake Bay observations

212 Model simulations were compared to cruise data from the CBP for 2004 and 2005 from 213 13 stations along the main stem of the Bay (Table 2, Fig. 2). The years 2004 and 2005 214 were selected to represent relatively wet and average years, respectively, and the 13 215 stations were chosen as they have been found to offer optimal estimates of Bay-wide 216 hypoxic volume (Bever et al., 2013). Stations were sampled on up to 34 cruises over the 217 two years (Table 2), generally twice a month from April to August and once a month for 218 the remainder of the year. Observational data can be downloaded from the CBP Water 219 Quality Database (http://www.chesapeakebay.net/data/downloads/cbp\_water\_quality\_ 220 database\_1984\_present). Variables downloaded from the CBP website and used in this 221 study were temperature, salinity, DO, nitrate + nitrite (hereafter abbreviated as "nitrate"), 222 and chlorophyll *a* (hereafter abbreviated as "chlorophyll"). For most cruises, observations 223 of temperature, salinity, and DO were made at roughly 1 m intervals throughout the water 224 column, whereas observations of chlorophyll and nitrate were generally made only at the 225 surface, bottom, and sometimes one or two mid-water column locations. For further 226 information on available water quality observations, please see USEPA (2012). While 227 these observations were publicly available for model assessment during calibration of all 228 of the models, they represent a very small subset of the 30 years of EPA observations 229 across over 100 Bay stations. The models compared here were calibrated based on access 230 to the larger data set and for conditions in the Bay in general, not specifically for the 13 231 stations and two years considered here.

#### 232 **2.3** Calculation of stratification and mixed layer depth

Stratification of the density and oxygen fields was examined to identify the maximum gradient of the pycnocline and oxycline as well as the depth of the top of the pycnocline and oxycline. In open ocean studies, the depth of the top of stratification is commonly referred to as the mixed layer depth (MLD), although this term is less frequently used in the estuarine literature. As the research presented here distinguishes between the depths of the top of the pycnocline and that of the oxycline, these will be referred to respectively as the density ( $\rho$ ) mixed layer depth (MLD<sub> $\rho$ </sub>) and the oxygen mixed layer depth (MLD<sub>O</sub>). Density was calculated via a classical density formula that is also utilized by the CBP for use in the Chesapeake Bay (Fofonoff and Millard, 1983; USEPA, 2004) and is a function of temperature and salinity.

243 The CBP defines the top and bottom of stratification in order to distinguish individual 244 designated use areas for water quality management purposes (USEPA, 2004). They 245 suggest that the top of the pycnocline be defined as the shallowest occurrence of a density gradient of 0.1 kg  $m^{-4}$  or greater as resolved by CBP profile observations, which are 246 247 typically spaced at 0.5 to 2 m depth intervals. If density gradients throughout the water column are less than 0.1 kg m<sup>-4</sup>, they define the water to be unstratified. The 0.1 kg m<sup>-4</sup> 248 249 threshold definition is designed to identify any initiation of stratification that may serve to 250 cut off vertical mixing from a nearly perfectly well mixed layer.

251 While the CBP definition described above delineates between designated use boundaries 252 according to density, our research focuses on the relationship between the pycnocline and 253 oxycline, requiring an alternate definition that can be applied to both the density and 254 oxygen distributions. In addition, the CBP definition often generates estimates for the 255 depth of the top of the pycnocline that are too shallow compared to the maximum depth 256 of surface mixing (Fig. 3). As a result, a percentage threshold criterion was developed 257 that identifies the bottom of the reasonably well-mixed layer, rather than perfectly mixed 258 layer, and is used in this analysis. The percentage threshold method defines a density or 259 DO profile as being stratified if a change of 10 % of the difference between the profile's 260 maximum and minimum values occurs within a single meter (Fig. 3). For example, if the 261 maximum DO concentration throughout the water column on an individual sampling date is 10 mg  $L^{-1}$  and the minimum concentration is 1 mg  $L^{-1}$ , stratification is defined to be 262 present if a difference of 0.9 mg  $L^{-1}$  is present within one meter. As recommended by the 263 264 CBP, the uppermost meter of the water column is not considered (USEPA, 2004). The 265 mixed layer depth is therefore defined as the shallowest level (below 1 m depth) where 266 stratification is identified. The minimum stratification criterion utilized in this analysis 267 requiring a profile to pass the 10% threshold also ensures that observations where very

little stratification exists do not bias the stratification results while also allowing for asingle criterion to be used across multiple stratification variables.

#### 270 2.4 Model skill metrics

271 Simulations of the Chesapeake Bay from the eight models described above were 272 statistically compared to historical monitoring data using a variety of skill metrics 273 including: root-mean squared difference (RMSD), bias, standard deviation, and 274 correlation coefficient. These metrics are illustrated on Taylor and target diagrams 275 (Taylor, 2001; Hofmann et al., 2008; Jolliff et al., 2009), which offer a compact way of 276 assessing model skill by displaying a number of different skill metrics. Target diagrams 277 illustrate the bias and total RMSD of model output, which Taylor diagrams do not. Taylor 278 diagrams include quantitative information on the standard deviations and correlations 279 between the model output and the observations, which target diagrams do not. Both 280 diagrams, however, represent unbiased RMSD, sometimes called "centered-pattern 281 RMSD". On target diagrams, a model symbol above the horizontal axis overestimates the 282 mean of the observations and a model symbol to the right of the vertical axis 283 overestimates the variability of the observations. (See Hofmann et al. (2008) and Jolliff et 284 al. (2009) for a more detailed description of these diagrams.) On Taylor diagrams, a 285 model symbol lying on the horizontal axis exactly correlates to the observations and a 286 model symbol further from the origin than the observation symbol overestimates the 287 standard deviation of the observations. (See Taylor et al. (2001) for a more detailed 288 description of these diagrams).

289 Taylor and target diagrams presented here are normalized to the standard deviation of the 290 observations, allowing multiple variables be represented on the same plot. This also 291 conveniently allows the unit circle on a target diagram to represent the skill of a model 292 defined as the mean of the observations. In effect, this means that if a model falls within 293 the unit circle, it exhibits a skill that is greater than the skill obtained if one were to 294 simply use the mean of the observations. The Taylor and target plots are either temporal 295 (displaying model skill at a single station over the study period) or spatial (displaying 296 model skill during a single month over the entire set of study stations). In addition,

summary diagrams are presented which combine both temporal (examining the seasonal
changes at each individual station) and spatial (examining differences across the Bay
during an individual month) variability.

300 Model skill was assessed using the hourly model output (daily for CH3D-ICM 301 chlorophyll and nitrate) that was nearest in time to that of the observation and from the 302 grid cell that encompassed the observation location. For months with two observations, 303 each observation was individually matched to the model output and the skill statistics 304 from those comparisons were averaged for that month. The native horizontal resolution 305 and bathymetry of the individual model grids was preserved in the comparison so as not 306 to bias the analysis through varying interpolation methodologies. For stratification 307 variables, the models and observations were interpolated to a 1m vertical grid that 308 extended only as deep as the individual models' bathymetry or deepest observation in 309 order to preserve the differences in bathymetric grids while allowing for a direct 310 comparison of the observations to the models. Model-data comparisons at the bottom of 311 the water column were not necessarily based on the same depths, since in many cases the 312 modeled bathymetry was shallower (or at times, deeper) than the deepest data point at a 313 given station. In order to avoid issues with extrapolation and/or grid stretching, data at the 314 bottom of the water column were always compared with model estimates from the 315 deepest grid cell provided by each particular model. Model-data comparisons for 316 stratification and mixed layer depths only included stations and times for which 317 stratification was defined to exist in both the observed and simulated fields.

#### 318 3 Results

An analysis of model skill of the combined temporal and spatial variability of DO at the surface and bottom of the water column, as well as at the observed  $MLD_O$ , indicates that all models, regardless of biogeochemical complexity or spatial resolution, exhibit a high degree of skill in reproducing observed DO (Fig. 4). Specifically, all models produce DO concentrations at the surface and bottom that have a normalized total RMSD less than one. The same is true for nearly all models for DO at the observed  $MLD_O$ . However, most models underestimate observed DO both at the surface and at the  $MLD_O$  (Fig. 4a). 326 The correlation between the observed and modeled DO is relatively constant with depth 327 (Fig. 4b), though on average slightly higher at the bottom (0.85) than at the surface 328 (0.80). Further, on average, the models simulate DO at the surface and bottom better than 329 they do at the  $MLD_{\Omega}$ . No statistical difference exists between the skill of models that 330 utilize a full biogeochemical component and those that utilize the simple constant-331 respiration oxygen parameterization. Based on an analysis of variance (ANOVA) 332 comparing the full biogeochemical models to the CRM models, the two model types do 333 not perform differently in terms of their ability to reproduce the combined temporal and 334 spatial variability of bottom DO as measured by total RMSD (p = 0.48). Overall, Model 335 M (the mean of the 8 models) consistently performs better than any individual model

across all depths examined (Fig. 4).

337 The monthly temporal variability of bottom DO at each station over the two years studied 338 is resolved similarly well by all of the models (Fig. 5a), but the models have difficulty 339 simulating spatial DO variability during each month (Fig. 5b). Due to the stations chosen 340 for this analysis (Fig. 2), the spatial variability being examined here is essentially the 341 north to south variability. Most models exhibit a latitudinal gradient with respect to their 342 skill in reproducing the temporal variability of bottom DO, with models overestimating 343 DO at the more northern stations (Fig. 5a). Some models differ in their ability to 344 reproduce summer (May to September) DO concentrations and winter (October to April) 345 DO concentrations (Fig. 5b). Models B, F, and G all distinctively overestimate mean DO 346 in the summer compared to the winter. In contrast, Models A and C perform similarly 347 well in both seasons (Fig. 5b). In addition, all three constant respiration models as well as 348 Models D and E substantially underestimate DO at several stations in the winter.

All eight models generally resolve the pycnocline and oxycline with similar skill (Fig. 6).

350 All models consistently underestimate the mean and standard deviation of the maximum

351 strength of stratification within the pycnocline and oxycline, defined herein as the

352 maximum vertical gradients of density and oxygen (Fig. 6a). All models, except for

353 Model A (see Sect. 4.2), also underestimate the mixed layer depth, regardless of whether

it is computed in terms of density or oxygen. (Note that these model symbols in Fig. 6a

are located above the y axis despite this negative bias in MLD because the vertical

coordinate system is oriented upwards.) Thus the models are producing stratification that
is both weaker than observed and higher (shallower) in the water column. The correlation
coefficient for these metrics is low, ranging between 0.1–0.6, and indicates that all
models are missing the majority of variability associated with the magnitude and location
of the pycnocline and oxycline (Fig. 6b). However, there is slightly more consistency and
better correlation coefficients among the models for the strength of stratification than the
depth of the mixed layers.

363 All eight models are also characterized by similar skill in representing the temporal and

364 spatial variability of density stratification and  $MLD_{\rho}$  (Fig. 7). There is a latitudinal

365 difference in skill of the models in reproducing the magnitude of the pycnocline and

366  $MLD_{\rho}$ , with model skill generally lower at the northern stations (Fig. 7a). Contrary to the

367 pattern shown for bottom DO (Fig. 5b), none of the models exhibit a significant seasonal

368 pattern between summer and winter in reproducing spatial variability of  $d\rho/dz$  or MLD<sub> $\rho$ </sub>

369 (Fig. 7b). However, Model A differentiates itself from the rest of the models in its pattern

of skill at reproducing the spatial and temporal variability of the  $MLD_{\rho}$  (see Sect. 4.2).

371 Temporal and spatial patterns for oxycline stratification (dO/dz) and  $MLD_O$  closely match

372 those of  $d\rho/dz$  and MLD<sub> $\rho$ </sub> (not shown).

373 All eight models reproduce the variability of bottom DO better than the variables that are 374 generally thought of as being the primary drivers of hypoxic conditions, including 375 stratification (Fig. 6), salinity, chlorophyll and nitrate (Fig. 8, Table 3). However, all 376 models reproduce patterns in temperature across the Bay and through time better than any 377 of the other variables in this model comparison (Fig. 8). All eight models, as well as the 378 Model Mean, are characterized by very low bias in modeled temperature, and correlation 379 coefficients of approximately 0.99; this high skill results from the very strong and 380 predictable seasonal temperature variability. Even though the five models with full 381 biogeochemical components (Models A, B, C, D, and E) are characterized by large 382 differences in their mechanistic approaches to modeling nitrate and chlorophyll, they 383 produce similar total RMSDs for all of the variables examined at both the surface and at 384 the bottom (Table 3).

385 The mean of the eight models (Model M) has a higher model skill (lower RMSD) than 386 any individual model across nearly every variable examined (Table 3). In addition, for 387 nearly all observations at all stations, the 95 % confidence interval of all model hindcasts 388 encapsulates the observed bottom DO concentration (Fig. 9), even though any individual 389 model may overestimate or underestimate observed DO. Models generally fall into 390 greater agreement during the summer, when DO is low, and into lesser agreement in the 391 winter when DO is replete. While this study does not allow for a true interannual 392 comparison, it is interesting that at station CB4.1C whereas the model ensemble closely 393 matches the timing of the drawdown of DO in the spring of 2004 (Fig. 9), it produces a 394 summer rather than spring initiation of hypoxic conditions in 2005. In addition, the model 395 ensemble produces a premature relaxing of hypoxic conditions for both years at this 396 observation station.

397 In order to better understand the impact of stratification on DO concentrations throughout 398 the water column, the relationship between the observed pycnocline strength and  $MLD_0$ 399 were compared to the observed oxycline strength and  $MLD_{\Omega}$ . Observations from 1998 to 400 2006 demonstrate that while there is not a strong correlation between the strengths of the 401 pycnocline and oxycline, there is a very strong correlation between MLD<sub>o</sub> and MLD<sub>O</sub> 402 (Fig. 10). Depending on the criteria used for defining the existence of stratification (see Sect. 2.3), the correlation of the pycnocline and oxycline strengths range between  $r^2 =$ 403 0.18 to 0.26 and the correlations of MLD<sub> $\rho$ </sub> and MLD<sub>O</sub> range between  $r^2 = 0.51$  to 0.82 404 405 (Table 4). Furthermore, correlation of the relationship between the  $MLD_{\rho}$  and  $MLD_{O}$  is 406 stronger for more severe stratification (Table 4). The relationship between the two mixed 407 layer depths is biased towards the MLD<sub>O</sub> being slightly located deeper in the water 408 column than the MLD<sub>0</sub>. As the cut-off criteria for the existence of stratification becomes 409 more stringent, the relationship becomes closer to 1:1.

#### 410 4 Discussion

#### 411 **4.1** How does the skill of various hydrodynamically-based DO models compare?

412 – In examining the eight 3-D models in this study, there is not a statistical

difference between the ability of simple and complex models to simulate the mean
and monthly variability of bottom DO; in addition, models with higher spatial
resolution do not necessarily produce better estimates of DO.

416 Models currently simulating hypoxia throughout Chesapeake Bay compute oxygen 417 concentrations in essentially two distinct ways: they either utilize a simple constant 418 respiration model or a full biogeochemical model. In this study, the relative skill of both 419 types of models is compared. Specifically, in examining results of the comparison 420 between five biogeochemical models (A, B, C, D, and E) and three simplistic constant 421 respiration models (F, G, and H), the two groups of models performed statistically similar 422 in their skill of reproducing bottom DO concentrations (Fig. 3, Table 3). These results 423 support those of Bever et al. (2013) who compared three constant respiration models with 424 the CBP regulatory model (Model A) and similarly found that all four of the models were 425 equally skillful in terms of reproducing the seasonal variability in bottom DO throughout 426 the Bay in 2004 and 2005. Consistent with the results of Scully (2013), this result implies 427 that the seasonal variability of DO in the Chesapeake Bay is primarily dependent on 428 underlying hydrodynamic mechanisms which are nearly identical for all eight models, 429 rather than on aspects related to the biogeochemical cycling which vary dramatically 430 between models and in fact are constant in three of the eight models. It should be noted, 431 however, that the two years studied here were relatively wet years and an analysis of dry 432 years may offer different results.

433 Many previous studies have examined the costs and benefits of adding complexity to 434 biogeochemical models. For example, increasing biogeochemical complexity has been 435 found to improve skill in some biogeochemical data assimilative parameter optimization 436 studies (Friedrichs et al., 2006, 2007; Lehmann et al., 2009; Bagniewski et al., 2011; 437 Ward et al., 2013; Xiao and Friedrichs, 2014). The additional parameters associated with 438 increased complexity generally provide more parameters that are available for additional 439 tuning and subsequent improved model-data agreement. This is in contrast to the results 440 of this analysis demonstrating that increased biogeochemical complexity does not 441 necessarily improve model-data agreement. In this case the increase in model complexity 442 has likely outpaced the ability of the researchers to fully tune the model to the available

443 observations. However, even past studies that have invoked formal parameter 444 optimization methodologies such as genetic algorithms and variational adjoint methods 445 (Friedrichs et al., 2007; Ward et al., 2010; Xiao and Friedrichs, 2014) have found that 446 under certain conditions, adding too much complexity does not necessarily improve 447 model skill and in fact can decrease model skill and portability, primarily due to artifacts 448 resulting from overtuning. This mirrors findings from the larger ecosystem modeling 449 community where the best-fit models are often those with intermediate complexity 450 (Fulton et al., 2003).

451 In this study, horizontal grid resolution differed significantly between model 452 implementations, with the most highly resolved grid (Model G) including more than nine 453 times more grid cells than the lower resolution grids (Table 1). A certain degree of 454 resolution is clearly required to successfully simulate dynamic processes, and a model 455 with 8–10 km resolution will not be able to correctly simulate the hydrodynamic 456 processes within the Bay (Feng et al., 2015). However, an increase in horizontal grid 457 resolution from  $\sim 1.8$  to  $\sim 0.6$  km, which results in a run-time change of a factor of nine, 458 or possibly of 27 if the time step is accordingly decreased by a factor of three, does not 459 necessarily result in a significant improvement in simulation skill of either stratification 460 or bottom oxygen. Although not shown here, additional sensitivity experiments with 461 Model G revealed that doubling the vertical resolution of this model had no significant 462 effect on the model's ability to resolve the depth of stratification or the maximum 463 magnitude of stratification. Thus, when selecting the optimal model resolution for a 464 simulation, it is critical to weigh the advantages of increased resolution with the increased 465 time required for simulation. With a given level of computational resources, fewer 466 sensitivity experiments can be conducted with a model using a more highly resolved grid.

467 Accurately simulating the observed spatial variability of DO (Fig. 4b) was a greater

468 challenge than simulating the temporal variability of DO (Fig. 4a) for all eight models

469 participating in this intercomparison. This is especially true in the winter months when

470 the vast majority of the Bay is oxygen replete and the models have difficulty representing

471 the observed variability from station to station. The majority of the models tend to

472 slightly overestimate mean bottom DO in the summer whereas multiple models (e.g.,

Models D, E, F, and G) exhibit a strong negative bias during January and/or February of
2005, primarily at stations in the middle to southern portion of the Bay's deep channel.
Interestingly, increased biological complexity and higher grid resolution do not
completely resolve this issue, as this is true for models utilizing full biogeochemical
models (Models D, E) as well as those using highly resolved model grids (Model G). This
is likely due to the ephemeral nature of the biological divers of DO.

479 The strong performance of the constant respiration models implies that these models may 480 be excellent candidates for providing short-term bottom oxygen forecasts. The high DO 481 skill of the CRM models primarily results from the fact that seasonal variations in 482 physical processes (primarily wind mixing and temperature) play a dominant role in 483 controlling the seasonal cycle of oxygen (Scully, 2013). Because the underlying 484 hydrodynamic models all use similar physical forcing, the constant respiration models are 485 able to simulate the seasonal cycle of DO with similar skill as the more complex 486 biogeochemical models. As a result, these simple models that are easier to tune and 487 require less in the way of computational resources than full biogeochemical models, may 488 be efficiently used to produce short-term (on the order of days) DO forecasts. On the 489 contrary, the more complex full biogeochemical models will be necessary for scenario-490 based and long-term (on the order of months to years) forecasting which requires that 491 models respond to prescribed changes in the biogeochemical environment, such as 492 increased rates of nutrient loading due to changes in land use, land cover, and/or climate.

# 493 4.2 How does model skill of DO compare to that of the primary drivers of DO494 variability?

- 495 Overall, model DO skill is greater than that of the variables generally considered
- to drive DO variability, such as stratification, salinity, mixed layer depth,
- 497 chlorophyll, and nitrate; only modeled temperature has higher skill than modeled498 DO.

Since dissolved oxygen concentrations in the Chesapeake Bay are controlled by physical
processes (e.g., advection, wind mixing, heating/cooling, and stratification), as well as
biological processes (e.g., photosynthesis and respiration), it is critical to understand the

502 skill of the models in terms of how well they reproduce the many factors influencing 503 oxygen concentrations. As expected, the five models containing a specific 504 biogeochemical model component had more difficulty simulating the observed 505 chlorophyll and nitrate concentrations than the physical variables (temperature and 506 salinity), both at the surface (Table 3) and the bottom (Fig. 8). Replicating the correct 507 location, magnitude, and timing of phytoplankton blooms and nutrient cycling is a 508 complex issue, and as a result, these features are generally not well simulated in the 509 models. While the models generally simulate the total amount of chlorophyll adequately, 510 it is more uniformly spatially distributed in the models rather than in patchy blooms as in 511 nature, leading to the underestimation of chlorophyll variability across all models. 512 Although all models produced a relatively high correlation between observed and 513 modeled temperature and salinity (Fig. 8), the correlation coefficients for chlorophyll and 514 nitrate were much lower. The correlations for observed vs. modeled DO was more similar 515 to that of the physical variables (temperature, salinity) than the biological variables 516 (chlorophyll and nitrate), highlighting that the seasonal variability in bottom DO is 517 regulated more by physical than biological factors. This also explains the success of the 518 constant respiration models, which by definition contain no biological variability yet 519 reproduce DO variability nearly as well as the most complex biogeochemical models.

520 In this study, model skill was also considerably higher for bottom oxygen than it was for 521 the vertical gradient of stratification and mixed layer depths (Figs. 6 and 8). The 522 underestimation of the vertical gradient across all models is largely due to the numerical 523 diffusion that characterizes all of these hydrodynamic models, but may also be partially 524 due to an underestimation of the winds or a lack of diffuse freshwater input around the 525 Bay. Even though the models all underestimated the strength of stratification (Figs. 4 and 526 6), modeled stratification in summer was strong enough to prevent mixing with the 527 relatively well-oxygenated surface waters. This result suggests, somewhat surprisingly, 528 that simulating the correct vertical gradient of stratification is not absolutely necessary for 529 skillful bottom DO simulations. Models need only simulate *enough* stratification to 530 effectively cut off vertical mixing in order to develop an isolated bottom layer that can 531 then experience a draw down in oxygen via respiration. In addition, the models must also

correctly simulate the horizontal advection of oxygen (Scully, 2013; Li et al., 2015). The
fact that bottom DO is simulated so well by the eight models analyzed here suggests that
not only is the advection of oxygen well represented in the models, but also the strength
of stratification, i.e., the maximum vertical gradients of density and oxygen, produced by
these models is sufficient. Thus, although novel and somewhat unexpected, these results
are not contradictory to previous studies demonstrating the importance stratification plays
in initiating summer hypoxia in the Chesapeake Bay (Murphy et al., 2011).

539 Model skill in terms of reproducing observed mixed layer depths was likewise much

be lower than model skill of reproducing observed oxygen concentrations. All models,

541 except Model A, produced mixed layer depths ( $MLD_0$  and  $MLD_0$ ) that were generally

too shallow in the water column (Fig. 6a). Note that Model A is a regulatory model that

has been used for many years by the Chesapeake Bay Program, and has thus undergone

544 more extensive calibration aimed at matching the mean salinity and oxygen

545 characteristics of the Bay (Cerco and Cole, 1993). Furthermore, Model A employs a z

546 grid that matches the bathymetry in trench areas better than the sigma grids used by the

547 other models. Although Model A produced mixed layer depths that were generally in the

548 correct location within the water column (Fig. 6a), they were too variable (Fig. 6b). This

variability may partly be a result of the 1.5m *z* grid employed by Model A causing large

550 jumps between vertical grid cells and hence resulting in overestimates of MLD

variability. All other models use sigma grids typically with more highly resolved vertical

resolution at the depth of maximum stratification.

553 The two variables for which the models have greatest skill are DO and temperature (Fig.

8). This is because oxygen variability is driven primarily by seasonal variability in

physical processes such as solubility and wind mixing and to a lesser degree by

variability in oxygen consumption (Scully, 2013). As a result, the models using a

557 constant mean respiration rate produce as realistic hypoxia simulations as the

558 biogeochemically complex models. Observations clearly show this strong seasonal

variability in bottom DO (Fig. 11a) and, to a slightly lesser extent, clear seasonal

variability in DO at the bottom of the bottom of the oxygen mixed layer ( $MLD_0$ ; Fig.

561 11b). But a seasonal cycle is not manifested in the MLD<sub>0</sub> itself (Fig. 11c). The lack of

such a strong seasonal cycle in the observed mixed layer depths makes this a more

563 difficult variable for the models to simulate. As a result, the models can relatively

skillfully simulate the combined spatial and temporal variability of DO while

565 simultaneously missing the  $MLD_0$ .

#### 566 **4.3** Why is it important for DO models to simulate the MLD<sub>O</sub> correctly?

Most of the aerobic habitat in the Bay during the summer is located above the
 MLD<sub>O</sub>, thus it is critical for living resource managers to use models that accurately
 simulate this variable.

570 On average, the models miss the observed depth of the MLD<sub> $\Omega$ </sub> by 3.4m, which equates to 571 roughly a 60 % error in the modeled mixed layer depths. While the models have 572 difficulty simulating the MLD<sub>O</sub> throughout the entire year (Figs. 6 and 7b), the summer 573 months are when the mismatch has the greatest potential to impact the available habitat 574 for oxygen-dependent species. Each year during this time period low-oxygen waters 575 occupy nearly the entire water column below the mixed layer. At Station CB4.1C, a representative mesohaline deep trough station, the contours of low-oxygen  $(5mgL^{-1})$  and 576 hypoxic  $(2mgL^{-1})$  waters are located just below the MLD<sub>0</sub> from late spring until late fall 577 578 (Fig. 12). The severe depletion of oxygen below the mixed layer compresses the habitable 579 space at this station to roughly 10 m (from a maximum of 32 m) during the annual low-580 oxygen event.

581 The impact of habitat compression can be substantial, as many Bay species require DO 582 concentrations well above the traditional hypoxic threshold (USEPA, 2010). While not all 583 of the main stem stations develop hypoxic water each year, most mesohaline stations 584 experience a dramatic drawdown of oxygen to levels during the summer that effectively 585 remove a large portion of the Bay from habitable space (Murphy et al., 2011; Schlenger 586 et al., 2013). Studies have shown that some species modify their behavior based on the 587 oxycline depth, which acts to constrict the habitable space in the water column (Prince 588 and Goodyear, 2006; Pierson et al., 2009; Elliot et al., 2013). Since species can be

negatively impacted by low-DO concentrations as high as  $5mgL^{-1}$  (Breitburg, 2002;

590 Vaquer-Sunyer and Duarte, 2008; USEPA, 2010), the location of the oxycline is not only

591 important for habitat compression in the summer months, but can also be important in the

592 winter months when an occasional lack of vertical mixing can substantially decrease

593 bottom DO concentrations. Furthermore, in order to accurately estimate hypoxic volume,

- 594 models must correctly simulate the depth of the mixed layer, since the MLD<sub>O</sub> closely
- 595 follows the depth of the 2 mg  $L^{-1}$  contour.

# 4.4 How can DO simulations in the Bay be improved for management of waterquality and living resources?

# 598 – To better simulate DO conditions and summer habitat compression due to low-599 DO water, simulations of the depth of the top of the pycnocline (MLD<sub> $\rho$ </sub>) must be 600 improved.

601 Although the suite of models examined reproduce DO concentrations relatively well 602 overall (Fig. 4), the models typically overestimate summer habitat compression by 603 producing low DO concentrations too high in the water column (Fig. 6). Observations 604 from the Chesapeake Bay Program show a strong correlation between the depths of the 605 oxygen and density-defined mixed layers (Fig. 10b). The models analyzed here also 606 clearly exhibit a close relationship between their skill in simulating the depths of the 607 oxygen and density-defined mixed layers (Fig. 6). These strong relationships between the 608 depths of the oxygen and density-defined mixed layers result from the fact that the 609 pycnocline represents the physical barrier that leads to the development of the oxycline. 610 Therefore, the inability of the models to accurately simulate habitat compression is an 611 artifact of their lack of skill in simulating the depth of the density-defined mixed layer. In 612 contrast, the strength of density stratification is not well correlated to the strength of 613 oxygen stratification. This is because a relative wide range of intensities of density 614 stratification is still sufficient to cut off vertical mixing, leading to the observed draw-615 down in bottom DO. Thus, even though all models underestimate the strength of the 616 pycnocline, they still produce enough stratification to greatly reduce mixing. The results 617 from this paper thus indicate that to further improve DO simulations and better estimate

summertime habitat compression, it is even more critical for models to accurately
simulate the depth of the top of the pycnocline than to accurately simulate the absolute
strength of the pycnocline.

#### 621 **4.5** What is the utility of the multi-model ensemble and Model Mean?

622 - The multi-model ensemble approach allows for the development of a Model
623 Mean, which taken as its own model, is the most skilled model when examining the
624 combined suite of variables analyzed in this study.

625 The model skill assessment presented here demonstrates that the average of all eight 626 models, or five models in the case of chlorophyll and nitrate, does better than any 627 individual model if looking across the suite of variables analyzed. This finding is similar 628 to that of other studies that examined the value of the model mean from a multi-model 629 ensemble (e.g., Gneiting and Raftery, 2005; Hagedorn et al., 2005). While the concept of 630 using a multi-model ensemble has been most extensively employed by atmospheric, 631 climatic, and global circulation modelers, such as the Intergovernmental Panel on Climate 632 Change (e.g., Collins et al., 2013), the tool's utility for aquatic ecosystem modeling is 633 gaining traction (Meier et al., 2012; Trolle et al., 2014; Janssen et al., 2015). As models 634 are increasingly used in regulatory decisions regarding aquatic ecosystems, a cohort of 635 similarly skilled models can be used to help inform a set of confidence bounds around an 636 environmental forecast. Due to the restrictions placed on models used in regulatory 637 actions, utilization of a multi-model ensemble may not be realistic for all environmental 638 and resource managers; however, multiple models can be integrated into the decision-639 making process even when the ultimate decision must be based on a single model. For 640 example, a confidence interval plot could help identify where regulatory model output 641 might be acting out of sync with other skilled water quality models of the same system, 642 thereby informing managers of the potential shortfalls associated with the regulatory 643 model. Furthermore, if the models tend to be predicting similar DO concentrations, a 644 cohort of models could enhance the confidence in regulatory decisions based on a single 645 regulatory model (Friedrichs et al., 2012; Weller et al., 2013). Comparing multiple 646 models can also help inform how to better improve models in the future, as this study has

647 aimed to do.

#### 648 **5** Conclusions

649 All models analyzed here exhibited a high degree of skill in simulating dissolved oxygen 650 concentrations within the main stem of the Chesapeake Bay in two years corresponding 651 to relatively wet and average years. Their high skill results from the fact that physical 652 processes (e.g., solubility, wind-mixing, and advection) exert a first order influence on 653 the seasonal cycle of oxygen. As a result, the models' ability to reproduce dissolved 654 oxygen concentrations is independent of the complexity of the biogeochemical 655 parameterizations: the simplest constant respiration models were found to reproduce 656 observed oxygen concentrations as well as the most biologically complex models. 657 Essentially, all models are equally capable of respiring most of the available oxygen in 658 the lower water column during summer.

659 This study also suggests that for use as management tools for water quality and living 660 resources, it is more critical for these models to adequately resolve the depth of the mixed 661 layer than the absolute strength of stratification (as long as modeled stratification is 662 strong enough to limit vertical mixing). This is critical because observations show that 663 during warmer months, oxygen-depleted water fills the water column to where 664 stratification limits further mixing, which effectively cuts off waters below the mixed 665 layer for use by the majority of the Chesapeake Bay's most recognized and valued living 666 resources. These results furthermore suggest that modelers should focus their efforts on 667 improving the hydrodynamics of their models in an effort to improve simulations of 668 mixed layer depth dynamics and variability.

669 These findings have significant ramifications for short-term bottom DO forecasts, which

670 may be successful with very simple oxygen parameterizations embedded in

671 hydrodynamic models. In contrast, scenario-based water quality forecasts are likely to

benefit from more complex models, which must adequately reproduce the longer-term

673 response of the oxygen field to changes in nutrient and organic matter loads. This study

also helps to demonstrate how multiple community models from governmental agencies

and academic institutions may be used together to provide a model mean and a set of

- 676 confidence bounds for regulatory model results that could be used to inform management
- 677 decisions.
- 678
- 679

680 681 682 683	<i>Acknowledgements</i> . This work was supported by the NOAA IOOS program as part of the Coastal Ocean Modeling Testbed. We thank Yun Li and Younjoo Lee for help with the ROMS-RCA simulations used in this analysis and Ray Najjar for his insights and comments. This is VIMS contribution 3520 and UMCES contribution 5130.
684	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
696	

#### 697 **References**

- Bagniewski, W., Fennel, K., Perry, M. J., and D'Asaro, E. A.: Optimizing models of the North
  Atlantic spring bloom using physical, chemical and bio-optical observations from a Lagrangian
  float, Biogeosciences, 8, 1291–1307, doi:10.5194/bg-8-1291-2011, 2011.
- 701 Bever, A. J., Friedrichs, M. A. M., Friedrichs, C. T., Scully, M. E., and Lanerolle, L. W.:
- 702 Combining observations and numerical model results to improve estimates of hypoxic volume
- within the Chesapeake Bay, USA, J. Geophys. Res-Oceans, 118, 4924–4944,
- 704 doi:10.1002/jgrc.20331, 2013.
- Boesch, D. F., Brinsfield, R. B., and Magnien, R. E.: Chesapeake Bay Eutrophication: scientific
  understanding, ecosystem restoration, and challenges for agriculture, J. Environ. Qual., 30, 303–
  320, 2001.
- Breitburg, D.: Effects of hypoxia, and the balance between hypoxia and enrichment, on coastal
  fishes and fisheries, Estuaries, 25, 767–781, 2002.
- Breitburg, D. L., Loher, T., Pacey, C. A., and Gerstein, A.: Varying effects of low dissolved
  oxygen on trophic interactions in an estuarine food web, Ecol. Monogr., 67, 489–507, 1997.
- 712 Brown, C. W., Hood, R. R., Long, W., Jacobs, J., Ramers, D. L., Wazniak, C., Wiggert, J. D.,
- 713 Wood, R., and Xu, J.: Ecological forecasting in Chesapeake Bay: using a mechanistic-empirical
- 714 modeling approach, J. Marine Syst., 125, 113–125, doi:10.1016/j.jmarsys.2012.12.007, 2013.
- 715 Buchheister, A., Bonzek, C. F., Gartland, J., and Latour, R. J.: Patterns and drivers of the
- demersal fish community of Chesapeake Bay, Mar. Ecol.-Prog. Ser., 481, 161–180,
  doi:10.3354/meps10253, 2013.
- Cerco, C., Johnson, B., and Wang, H.: Tributary Refinements to the Chesapeake Bay Model,
   ERDC TR-02-4, US Army Engineer Research and Development Center, Vicksburg, MS, 2002.
- Cerco, C., Kim, S.-C., and Noel, M.: The 2010 Chesapeake Bay Eutrophication Model a Report
  to the US Environmental Protection Agency Chesapeake Bay Program and to The US Army
  Engineer Baltimore District, US Army Engineer Research and Development Center, Vicksburg,
- 723 MS, 2010.
- Cerco, C. F. and Cole, T.: Three-dimensional eutrophication model of Chesapeake Bay, J.
  Environ. Eng.-ASCE, 119, 1006–10025, 1993.
- Cerco, C. F. and Noel, M. R.: The 2002 Chesapeake Bay Eutrophication Model, EPA 903-R-04004, US Army Corps of Engineers, Waterways Experiment Stations, Vicksburg, MS, 2004.
- 728 Cerco, C. F. and Noel, M. R.: Twenty-one-year simulation of Chesapeake Bay water quality
- view one-year simulation of chesapeake bay water quarty
   using the CE-QUAL-ICM eutrophication model, J. Am. Water Resour. As., 49, 1119–1133,
   doi:10.1111/jawr.12107, 2013.
- 731 Collins, M., Knutti, R., Arblaster, J., Dufresne, J.-L., Fichefet, T., Friedlingstein, P., Gao, X.,
- 732 Gutowski, W. J., Johns, T., Krinner, G., Shongwe, M., Tebaldi, C., Weaver, A. J., and Wehner,
- 733 M.: Long-term climate change: projections, commitments and irreversibility, in: Climate Change

- 734 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment
- Report of the Intergovernmental Panel on Climate Change, Cambridge University Press,
- 736 Cambridge, United Kingdom and New York, NY, USA, 1029–1136, 2013.
- Cooper, S. R. and Brush, G. S.: Long-term history of Chesapeake Bay anoxia, Science, 254, 992–
  996, 1991.
- Cooper, S. R. and Brush, G. S.: A 2,500-year history of anoxia and eutrophication in Chesapeake
  Bay, Estuaries, 16, 617–626, 1993.
- 741 Diaz, R. J.: Overview of hypoxia around the world, J. Environ. Qual., 30, 275–281, 2001.
- Diaz. R. J. and Rosenberg, R.: Spreading dead zones and consequences for marine ecosystems,
   Science, 321, 926–929, doi:10.1126/science.1156401, 2008.
- Du, J. and Shen, J.: Decoupling the influence of biological and physical processes on the
- dissolved oxygen in the Chesapeake Bay, J. Geophys. Res.-Oceans, 120, 78–93,
  doi:10.1002/2014JC010422, 2015.
- , io administration (2011) (2010) (22, 2010)
- Ekau, W., Auel, H., Pörtner, H.-O., and Gilbert, D.: Impacts of hypoxia on the structure and
  processes in pelagic communities (zooplankton, macro-invertebrates and fish), Biogeosciences, 7,
  1669–1699, doi:10.5194/bg-7-1669-2010, 2010.
- Elliott, D. T., Pierson, J. J., Roman, M. R.: Predicting the effects of coastal hypoxia on vital rates
  of the planktonic copepod *Acartia tonsa* dana, PLoS ONE, 8, e63987,
- 752 doi:10.1371/journal.pone.0063987, 2013.
- Feng, Y., Friedrichs, M. A. M., Wilkin, J., Tian, H., Yang, Q., Hofmann, E. E., Wiggert, J. D.,
- and Hood, R. R.: Chesapeake Bay nitrogen fluxes derived from a land-estuarine ocean
- biogeochemical modeling system: model description, evaluation, and nitrogen budgets, J.
  Geophys. Res.-Biogeo., 120, 1666–1695, doi:10.1002/2015JG002931, 2015.
- Fofonoff, N. P. and Millard, R. C.: Algorithms for Computations of Fundamental Properties of
  Seawater, UNESCO Technical Papers in Marine Science, 44, Paris, France, 53 pp., 1983.
- 759 Friedrichs, M., Sellner, K. G., and Johnston, M. A.: Using Multiple Models for Management in
- 760 the Chesapeake Bay: a Shallow Water Pilot Project, Chesapeake Bay Program Scientific and
- 761 Technical Advisory Committee Report, No. 12-003, Edgewater, MD, 2012.
- Friedrichs, M. A. M., Hood, R., and Wiggert, J.: Ecosystem model complexity versus physical
  forcing: quantification of their relative impact with assimilated Arabian Sea data, Deep-Sea Res.
  Pt. II, 53, 576–600, 2006.
- 701 11.11, 55, 570-000, 2000.
- Friedrichs, M. A. M., Dusenberry, J., Anderson, L., Armstrong, R., Chai, F., Christian, J., Doney,
- S. C., Dunne, J., Fujii, M., Hood, R., McGillicuddy, D., Moore, K., Schartau, M., Sptiz, Y. H.,
- and Wiggert, J.: Assessment of skill and portability in regional marine biogeochemical models:
- role of multiple phytoplankton groups, J. Geophys. Res., 112, C08001,
- 769 doi:10.1029/2006JC003852, 2007.
- Fulton, E. A., Smith, A. D. M., and Johnson, C. R.: Effect of complexity on marine ecosystem

- 771 models, Mar. Ecol.-Prog. Ser., 253, 1–16, 2003.
- Gilbert, D., Rabalais, N. N., Díaz, R. J., and Zhang, J.: Evidence for greater oxygen decline rates
  in the coastal ocean than in the open ocean, Biogeosciences, 7, 2283–2296, doi:10.5194/bg-72283-2010, 2010.
- Gneiting, T. and Raftery, A. E.: Weather forecasting with ensemble methods, Science, 310, 248–
  249, doi:10.1126/science.1115255, 2005.
- Hagedorn, R., Doblas-Reyes, F. J., and Palmer, T. N.: The rationale behind the success of multimodel ensembles in seasonal forecasting I. Basic concept, Tellus A, 57, 219–233,
  doi:10.1111/j.1600-0870.2005.00103.x, 2005.
- Hagy, J. D., Boyton, W. R., Keefe, C. W., and Wood, K. V.: Hypoxia in Chesapeake Bay, 1950–
  2001: long-term change in relation to nutrient loading and river flow, Estuaries, 27, 634–658,
  2004.
- Haidvogel, D. B., Arango, H., Budgell, W. P., Cornuelle, B. D., Curchitser, E., Di Lorenzo, E.,
- Fennel, K., Geyer, W. R., Hermann, A. J., Lanerolle, L., Levin, J., McWilliams, J. C., Miller, A.
- J., Moore, A. M., Powell, T. M., Shchepetkin, A. F., Sherwood, C. R., Signell, R. P., Warner, J.
- C., and Wilkin, J.: Ocean forecasting in terrain-following coordinates: formulation and skill
  assessment of the Regional Ocean Modeling System, J. Comput. Phys., 227, 3595–3624,
- 788 doi:10.1016/j.jcp.2007.06.016, 2008.
- Harding Jr., L. W. and Perry, E. S.: Long-term increase of phytoplankton biomass in Chesapeake
  Bay, 1950–1994, Mar. Ecol.-Prog. Ser., 157, 39–52, 1997.
- Harding Jr., L. W., Gallegos, C. L., Perry, E. S., Miller, W. D., Adolf, J. E., Mallonee, M. E., and
- Paerl, H. W.: Long-term trends of nutrients and phytoplankton in Chesapeake Bay, Estuaries
  Coasts, doi:10.1007/s12237-015-0023-7, 2015.
- Hofmann, E. E., Druon, J., Fennel, K., Friedrichs, M., Haidvogel, D., Lee, C., Mannino, A.,
- 795 McClain, C., Najjar, R., O'Reilly, J., Pollard, D., Previdi, M., Seitzinger, S., Siewert, J.,
- Signorini, S., and Wilkin, J.: Eastern US continental shelf carbon budget: integrating models, data
- assimilation, and analysis, Oceanography, 21, 86–104, doi:10.5670/oceanog.2008.70, 2008.
- Hong, B. and Shen, J.: Responses of estuarine salinity and transport processes to potential future
  sea-level rise in the Chesapeake Bay, Estuar. Coast. Shelf S., 104–105, 33–45,
- 800 doi:10.1016/j.ecss.2012.03.014, 2012.
- Hong, B. and Shen, J.: Linking dynamics of transport timescale and variations of hypoxia in the
   Chesapeake Bay, J. Geophys. Res.-Oceans, 118, 6017–6029, doi:10.1002/2013JC008859, 2013.
- Janssen, A. B. G., Arhonditsis, G. B., Beusen, A., Bolding, K., Bruce, L., Bruggeman, J.,
- 804 Couture, R.-M., Downing, A. S., Elliott, J. A., Frassl, M. A., Gal, G., Gerla, D. J., Hipsey, M. R.,
- Hu, F., Ives, S. C., Janse, J. J., Jeppsen, E., Johnk, K. D., Kneis, D., Kong, X., Kuiper, J. J.,
- Lehmann, M. K., Lemmen, C., Ozkundakci, D., Petzoldt, T., Rinke, K., Robson, B. J., Sachse, R.,
- 807 Schep, S. A., Schmid, M., Scholten, H., Teurlincx, S., Trolle, D., Troost, T. A., Van Dam, A. A.,
- 808 Van Gerven, L. P. A., Weijerman, M., Wells, S. A., and Mooij, W. M.: Exploring, exploiting and
- 809 evolving diversity of aquatic ecosystem models: a community perspective, Aquat. Ecol., 49, 513–

- 810 548, doi:10.1007/s10452-015-9544-1, 2015.
- Jiang, L. and Xia, M.: Dynamics of the Chesapeake Bay outflow plume: Realistic plume
- 812 simulations and its seasonal and interannual variability, J. Geophys. Res.-Oceans, 121,
- 813 doi:10.1002/2015JC011191, 2016.
- Jiang, L., Xia, M., Ludsin, S. A., Rutherford, E. S., Mason, D. M., Jarrin, J. M., and Pangle, K.
- 815 L.: Biophysical modeling assessment of the drivers for plankton dynamics in dressenid-colonized
- 816 western Lake Erie, Ecol. Model., 308, 18–33, 2015.
- Jolliff, J. K., Kindle, J. C., Schulman, I., Penta, B., Friedrichs, M. A. M., Helber, R., and Arnone,
- 818 R. A.: Summary diagrams for coupled hydrodynamic-ecosystem model skill assessment, J.
- 819 Marine Syst., 76, 64–82, doi:10.1016/j.jmarsys.2008.05.014, 2009.
- Keisman, J. and Shenk, G.: Total maximum daily load criteria assessment using monitoring and
  modeling data, J. Am. Water Resour. As., 49, 1134–1149, doi:10.1111/jawr.12111, 2013.
- 822 Keister, J. E., Houde, E. D., and Breitburg, D. L.: Effects of bottom-layer hypoxia on abundances
- and depth distributions of organisms in Patuxent River, Chesapeake Bay, Mar. Ecol.-Prog. Ser.,
- 824 205, 43–59, 2000.
- 825 Kemp, W. M., Boyton, W. R., Adolf, J. E., Boesch, D. F., Boicourt, W. C., Brush, G., Cornwell,
- J. C., Fisher, T. R., Gilbert, P. M., Hagy, J. D., Harding, L. W., Houde, E. D., Kimmel, D. G.,
- Miller, W. D., Newell, R. I. E., Roman, M. R., Smith, E. M., and Stevenson, J. C.: Eutrophication
  of Chesapeake Bay: historical trends and ecological interactions, Mar. Ecol.-Prog. Ser., 303, 1–
- 829 29, 2005.
- Kemp, W. M., Testa, J. M., Conley, D. J., Gilbert, D., and Hagy, J. D.: Temporal responses of
  coastal hypoxia to nutrient loading and physical controls, Biogeosciences, 6, 2985–3008,
  doi:10.5194/bg-6-2985-2009, 2009.
- Lanerolle, L. W., Patchen, R. C., and Aikman, F.: The Second Generation Chesapeake Bay
- 834 Operational Forecast System (CBOFS2): Model Development and Skill Assessment, TR- NOS-
- 835 CS-29, US Department of Commerce, National Oceanic and Atmospheric Administration,
  836 National Ocean Service, Office of Coast Survey, Coast Survey Development Laboratory, Silver
- 837 Spring, MD, 2011.
- Lehmann, M. K., Fennel, K., and He, R.: Statistical validation of a 3-D bio-physical model of the
  western North Atlantic, Biogeosciences, 6, 1961–1974, doi:10.5194/bg-6-1961-2009, 2009.
- 840 Levin, L. A., Ekau, W., Gooday, A. J., Jorissen, F., Middelburg, J. J., Naqvi, S. W. A., Neira, C.,
- Rabalais, N. N., and Zhang, J.: Effects of natural and human-induced hypoxia on coastal benthos,
  Biogeosciences, 6, 2063–2098, doi:10.5194/bg-6-2063-2009, 2009.
- Li, M., Zhong, L., and Boicourt, W. C.: Simulations of Chesapeake Bay estuary: sensitivity to
  turbulence mixing parameterizations and comparison with observations, J. Geophys. Res., 110,
  C12004, doi:10.1029/2004JC002585, 2005.
- Li, Y., Li, M., and Kemp, W. M.: A budget analysis of bottom-water dissolved oxygen in
  Chesapeake Bay, Estuar. Coast., 38, 2132–2148, doi:10.1007/s12237-014-9928-9, 2015.

- 848 Meier, H. E. M., Andersson, H. C., Arheimer, B., Blenckner, T., Chubarenko, B., Donnelly, C.,
- Eilola, K., Gustafsson, B. G., Hansson, A., Havenhand, J., Hoglund, A., Kuznetsov, I.,
- 850 MacKenzie, B. R., Muller-Karulis, B., Neumann, T., Niiranen, S., Piwowarczyk, J., Raudsepp,
- U., Reckermann, M., Ruoho-Airola, T., Savchuk, O. P., Schenk, F., Schimanke, A., Vali, G.,
- 852 Weslawski, J.-M., and Zorita, E.: Comparing reconstructed past variations and future projections
- 853 of the Baltic Sea ecosystem first results from multi-model ensemble simulations, Environ. Res.
- 854 Lett., 7, 034005, doi:10.1088/1748-9326/7/3/034005, 2012.
- 855 Meire, L., Soetaert, K. E. R., and Meysman, F. J. R.: Impact of global change on coastal oxygen
- dynamics and risk of hypoxia, Biogeosciences, 10, 2633–2653, doi:10.5194/bg-10-2633-2013,
  2013.
- 858 Murphy, R. R., Kemp, W. M., Ball, W. P.: Long-term trends in Chesapeake Bay seasonal
- hypoxia, stratification, and nutrient loading, Estuar. Coast., 34, 1293–1309, doi:10.1007/s12237011- 9413-7, 2011.
- 861 Najjar, R. G., Pyke, C. R., Adams, M. B., Breitburg, D., Hershner, C., Kemp, M., Howarth, R.,
- 862 Mulholland, M. R., Paolisso, M., Secor, D., Sellner, K., Wardrop, D., and Wood, R.: Potential
- 863 climate-change impacts on the Chesapeake Bay, Estuar. Coast. Shelf S., 86, 1–20,
- doi:10.1016/j.ecss.2009.09.026, 2010.
- Park, K., Kuo, A. Y., Shen, J., and Hamrick, J. M.: A three-dimensional Hydrodynamic
- Eutrophication Model (HEM-3D): description of water quality and sediment process submodels,in: Applied Marine Science and Ocean Engineering, Special Report, Virginia Institute of Marine
- 868 Science, Gloucester Point, VA, 327, 1995.
- Pierson, J. J., Roman, M. R., Kimmel, D. G., Boicourt, W. C., and Zhang, X. S.: Quantifying
- changes in the vertical distribution of mesozooplankton in response to hypoxic bottom waters, J.
  Exp. Mar. Biol. Ecol., 381, 74–79, 2009.
- Prince, E. D. and Goodyear, C. P.: Hypoxia-based habitat compression of tropical pelagic fishes,
  Fish. Oceanogr., 15, 451–464, doi:10.1111/j.1365-2419.2005.00393.x, 2006.
- 874 Riedel, B., Pados, T., Pretterebner, K., Schiemer, L., Steckbauer, A., Haselmair, A., Zuschin, M.,
- and Stachowitsch, M.: Effect of hypoxia and anoxia on invertebrate behaviour: ecological
- perspectives from species to community level, Biogeosciences, 11, 1491–1518, doi:10.5194/bg11-1491-2014, 2014.
- 878 Schlenger, A. J., North, E. W., Schlag, Z., Li, Y., Secor, D. H., Smith, K. A., and Niklitschek, E.
- 379 J.: Modeling the influence of hypoxia on the potential habitat of Atlantic sturgeon *Acipenser*
- 880 *oxyrinchus*: a comparison of two methods, Mar. Ecol.-Prog. Ser., 483, 257–272,
- doi:10.3354/meps10248, 2013.
- Scully, M. E.: The importance of climate variability to wind-driven modulation of hypoxia in
  Chesapeake Bay, J. Phys. Oceanogr., 40, 1435–1440, doi:10.1175/2010JPO4321.1, 2010.
- Scully, M. E.: Physical controls on hypoxia in Chesapeake Bay: a numerical modeling study, J.
  Geophys. Res.-Oceans, 118, 1239–1256, doi:10.1002/jgrc.20138, 2013.
- 886 Shenk, G. W. and Linker, L. C.: Development and application of the 2010 Chesapeake Bay

- 887 watershed total maximum daily load model, J. Am. Water Resour. As., 49, 1–15, 888 doi:10.1111/jawr.12109, 2013.
- 889 Shchepetkin, A. F. and McWilliams, J. C.: The Regional Ocean Modeling System (ROMS): a
- 890 split-explicit, free-surface, topography-following-coordinate oceanic model, Ocean Model., 9, 891 347-404, doi:10.1016/j.ocemod.2004.08.002, 2005.
- 892 Taylor, K. E.: Summarizing multiple aspects of models performance in a single diagram, J. 893 Geophys. Res., 106, 7183-7192, 2001.
- 894 Testa, J. M. and Kemp, W. M.: Spatial and temporal patterns of winter-spring oxygen depletion 895 in Chesapeake Bay bottom water, Estuar. Coast., 37, 1432–1448, doi:10.1007/s12237-014-9775-896 8,2014.
- 897 Testa, J. M., Li, Y., Lee, Y. J., Li, M., Brady, D. C., Di Toro, D. M., Kemp, W. M., and
- 898 Fitzpatrick, J. J.: Quantifying the effects of nutrient loading on dissolved O<sub>2</sub> cycling and hypoxia

899 in Chesapeake Bay using a coupled hydrodynamic-biogeochemical model, J. Marine Syst., 139,

- 900 139-158, doi:10.1016/j.jmarsys.2014.05.018, 2014.
- 901 Tian, H., Yang, Q., Najjar, R., Ren, W., Friedrichs, M. A. M., Hopkinson, C. S., and Pan, S.:
- 902 Anthropogenic and climatic influences on carbon fluxes from eastern North America to the
- 903 Atlantic Ocean: a process-based modeling study, J. Geophys. Res.-Biogeo., 120, 752–772,
- 904 doi:10.1002/2014JG002760, 2015.
- 905 Trolle, D., Elliott, J. A., Mooij, W. M., Janse, J. H., Bolding, K., Hamilton, D. P., and Jeppsen,
- 906 E.: Advancing projections of phytoplankton responses to climate change through ensemble

907 modeling, Environ. Modell. Softw., 61, 371-379, doi:10.1016/j.envsoft.2014.01.032, 2014.

908 USEPA: Ambient Water Quality Criteria for Dissolved Oxygen, Water Clarity, and Chlorophyll a

909 for the Chesapeake Bay and its Tidal Tributaries – 2004 Addendum, EPA 903-R-03-002, US

910 Environmental Protection Agency, USEPA Region III Chesapeake Bay Program Office,

- 911 Annapolis, MD, 2004.
- 912 USEPA: Chesapeake Bay Total Maximum Daily Load for Nitrogen, Phosphorus, and Sediment,
- 913 US Environmental Protection Agency, US Environmental Protection Agency Chesapeake Bay 914 Program Office, Annapolis, MD, 2010.
- 915 USEPA: Guide to Using Chesapeake Bay Program Water Quality Monitoring Data, EPA 903-R-
- 916 12-001, US Environmental Protection Agency, US Environmental Protection Agency Chesapeake 917 Bay Program, Annapolis, MD, 2012.
- 918 Vaquer-Sunyer, R. and Duarte, C. M.: Thresholds of hypoxia for marine biodiversity, P. Natl. 919 Acad. Sci. USA, 105, 15452–15457, doi:10.1073/pnas.0803833195, 2008.
- 920 Ward, B. A., Friedrichs, M. A. M., Anderson, T. R., and Oschlies, A.: Parameter optimization
- 921 techniques and the problem of underdetermination in marine biogeochemical models, J. Marine 922 Syst., 81, 34–43, doi:10.1016/j.jmarsys.2009.12.005, 2010.
- 923 Ward, B. A., Schartau, M., Oschlies, A., Martin, A. P., Follows, M. J., and Anderson, T. R.:
- 924 When is a biogeochemical model too complex? Objective model reduction and selection for

- 925 North Atlantic time-series sites, Prog. Oceanogr., 116, 49–65, doi:10.1016/j.pocean.2013.06.002,
  926 2013.
- 927 Weller, D., Benham, B., Friedrichs, M., Gardner, N., Hood, R., Najjar, R., Paolisso, M., Pasquale,
- 928 P., Sellner, K., and Shenk, G.: Multiple Models for Management in the Chesapeake Bay,
- 929 Chesapeake Bay Program Scientific and Technical Advisory Committee Workshop Report, No.
- 930 14-004, 25–26 February 2013.
- 931 Xiao, Y. and Friedrichs, M. A. M.: Using biogeochemical data assimilation to assess the relative
- 932 skill of multiple ecosystem models in the Mid-Atlantic Bight: effects of increasing the complexity
  933 of the planktonic food web, Biogeosciences, 11, 3015–3030, doi:10.5194/bg-11- 3015-2014,
  934 2014.
- Yu, J., Long, W., Wiggert, J. D., Lanerolle, L. W. J., Brown, C. W., Murtugudde, R., and Hood,
- R. R.: Climate forcing and salinity variability in Chesapeake Bay, USA, Estuar. Coast. Shelf S.,
- 937 35, 237–261, doi:10.1007/s12237-011-9423-5, 2012.
- 938 Yang, A., Tian, H., Friedrichs, M. A. M., Hopkinson, C. S, Lu, C., and Najjar, R. G.: Increased
- nitrogen export from eastern North America to the Atlantic Ocean due to climatic and
- 940 anthropogenic changes during 1901–2008, J. Geophys. Res.-Biogeo., 120, 1046-1068, doi:10.1002/2014/C002762.2015
- 941 doi:10.1002/2014JG002763, 2015.
- Yang, Q., Tian, H., Friedrichs, M. A. M., Liu, M., Li, X., and Yang, J.: Hydrological responses to
   climate and land-use changes along the North American east coast: a 110-year historical
- 944 reconstruction, J. Am. Water Resour. As., 51, 47–67, doi:10.1111/jawr.12232, 2015.
- 945 946 947 948 949 950 951 952 953 954 955 956 957 958 959 960 961 962 963 964
- 964
- 965 966

968 Table 1. Model characteristics.

Model	А	В	С	D	E	F	G	Н
Hydrodynamic model- DO model	CH3D- ICM	ChesROMS- ECB	ChesROMS- BGC	ROMS- RCA	FVCOM- ICM	ChesROMS- CRM	CBOFS- CRM	EFDC- CRM
Grid structure	Structured	Structured	Structured	Structured	Unstructured	Structured	Structured	Structured
Average wet-cell resolution	1 km	1.8 km	1.8 km	1.89 km	1.26 km	1.8 km	0.565 km	1.2 km
Vertical grid	1.52 m	20 sigma	20 sigma	20 sigma	10 sigma	20 sigma	20 sigma	20 sigma
River forcing	CBP Watershed Model	DLEM Watershed Model	USGS Data	USGS Data	USGS Data	USGS Data	USGS Data	USGS Data
Sub-tidal elevation at open boundary	Multiple efforts	Lewes, DE to Duck, NC	Lewes, DE to Duck, NC	Wachapreague, VA to Duck, NC	TPXO Tidal Model	Lewes, DE to Duck, NC	Ocean City, MD to Duck, NC	Lewes, DE to Duck, NC
Wind forcing	Multiple efforts	Thomas Point Light	NARR	NARR	NARR	NARR	NARR & NDBC buoys	NARR
Other atmospheric forcing	Multiple efforts	NARR	NARR	NARR	NARR	NARR	NARR	Norfolk & Baltimore Airports
Biogeochemical complexity	High; 5 phytoplk. groups	High; 1 phytoplk. group	High; 1 phytoplk. group	High; 2 phytoplk. groups	High; 3 phytoplk. groups	Low; constant respiration	Low; constant respiration	Low; constant respiration
Model citation	Cerco et al., 2010	Feng et al., 2015	Brown et al., 2013	Testa et al., 2014	Jiang and Xia, 2015	Scully, 2013	Lanerolle et al., 2011	Du and Shen, 2015
970								
971								
972								

Station	Latitude	Longitude	Station Depth	# of Cruises
CB3.2	39.1634 N	76.3063 W	12.1 m	34
CB3.3C	38.9951 N	76.3597 W	24.3 m	34
CB4.1C	38.8251 N	76.3997 W	32.3 m	34
CB4.2C	38.6448 N	76.4177 W	27.2 m	34
<b>CB4.3</b> C	38.5565 N	76.4347 W	26.9 m	34
<b>CB4.4</b>	38.4132 N	76.3430 W	30.3 m	34
CB5.1	38.3185 N	76.2930 W	34.1 m	34
CB5.2	38.13678N	76.2280 W	30.6 m	34
CB5.4	37.8001 N	76.1747 W	31.1 m	26
CB6.2	37.4868 N	76.1563 W	10.5 m	30
CB6.4	37.2365 N	76.2080 W	10.2 m	29
CB7.1	37.6835 N	75.9897 W	20.9 m	27
LE2.3	38.0215 N	76.3477 W	20.1 m	34

980 Table 2. Characteristics of observation stations (from USEPA, 2012).

, , <u>,</u>

## 993 Table 3. Mean and standard deviation (STD) of observations and total normalized RMSD

## for each model.

## 

					NT	1. 1.D.				
	Mean $\pm$		Normalized RMSD							
	STD of Obs	А	В	С	D	Е	F	G	Н	М
Surface Temp. (°C)	17.44±8.82	0.13	0.13	0.12	0.09	0.13	0.13	0.16	0.19	0.10
Bottom Temp. (°C)	$15.75 \pm 8.02$	0.24	0.35	0.35	0.23	0.22	0.35	0.17	0.19	0.23
Surface Salinity (PSU)	$10.92 \pm 4.32$	0.37	0.62	0.53	0.36	0.46	0.61	0.57	0.41	0.35
<b>Bottom Salinity</b> (PSU)	18.17±3.14	0.72	0.85	0.73	1.55	1.28	0.78	1.03	0.97	0.75
<b>Max.</b> $d\rho/dz$ (kg m <sup>-4</sup> )	~1.64±1.15	1.03	1.09	1.07	1.09	1.25	1.01	1.23	1.02	N/A
MLDp (m)	$\sim 5.32 \pm 3.99$	1.01	1.13	1.11	1.41	1.39	1.12	1.38	1.13	N/A
Surface DO $(mg L^{-1})$	9.74±2.15	0.67	0.58	0.89	0.80	1.00	0.63	0.64	0.69	0.57
<b>DO at MLD</b> <sub>O</sub> (mg $L^{-1}$ )	$\sim 8.44 \pm 2.53$	0.54	0.57	0.74	0.93	0.83	0.81	0.95	1.09	0.62
<b>Bottom DO</b> (mg $L^{-1}$ )	$4.42 \pm 3.61$	0.51	0.59	0.81	0.61	0.54	0.46	0.61	0.60	0.46
<b>Max. dDO/dz</b> (mg $L^{-1}$ m <sup>-1</sup> )	~1.81±1.12	1.19	1.21	1.34	1.09	1.35	1.12	1.23	1.19	N/A
MLDo (m)	~6.62±4.01	1.24	1.01	1.10	1.33	1.33	1.05	1.30	1.29	N/A
Surface Chl a (mg m <sup>-3</sup> )	$11.19 \pm 9.04$	0.92	1.22	1.60	1.23	0.89	N/A	N/A	N/A	1.16
Bottom Chl a (mg m <sup>-3</sup> )	9.02±11.52	0.87	1.10	1.07	1.05	1.01	N/A	N/A	N/A	0.90
<b>Surface Nitrate</b> (mmolN m <sup>-3</sup> )	$0.32 \pm 0.33$	0.61	0.79	1.03	0.61	0.52	N/A	N/A	N/A	0.79
<b>Bottom Nitrate</b> (mmolN m <sup>-3</sup> )	0.12±0.13	1.08	1.38	1.38	0.92	1.46	N/A	N/A	N/A	0.85
0.0.4										

999 Table 4. Pycnocline and oxycline correlation statistics (all correlations have p-values <<

### 1000 0.01).

Stratification Threshold Percentage	Max dρ/dz vs. Max dO/dz	MLDp vs. MLDo	Profiles with Stratification
10%	0.18	0.51	1613
15%	0.22	0.59	1303
20%	0.22	0.70	916
25%	0.26	0.82	575



1021 Figure 1. Map of the Chesapeake Bay and its watershed.



1034 Figure 2. Location of the CBP Water Quality Monitoring stations used in this study.



1046 Figure 3. Density and dissolved oxygen profiles for a mid-Bay station (CB4.1C) on (a)

1047 January 13, 2004 and (b) June 14, 2005, comparing the 0.1 kg  $m^{-4}$  stratification definition

1048 used by the CBP (MLD<sub>CBP</sub>) with the 10% threshold definitions used here for density

```
1049 (MLD_{\rho}) and oxygen (MLD_{O}).
```

1057



1058

1059 Figure 4. Normalized summary (a) target and (b) Taylor diagrams illustrating model skill

1060 of dissolved oxygen at the surface, MLD<sub>O</sub>, and bottom for 13 Chesapeake Bay stations in

1061 2004-2005. The "x" represents the skill of a model that perfectly reproduces the

1062 observations. The dotted, dashed-dot, and dashed lines on the Taylor diagram represent

1063 lines of constant standard deviation, correlation coefficient, and unbiased RMSD,

1064 respectively.





Figure 5. Normalized target diagrams for Models A-H demonstrating the (a) temporal and
(b) spatial skill in resolving the variability of bottom dissolved oxygen concentrations. In
(a) the individual dots represent the 13 stations along the main stem of the Chesapeake
Bay. In (b) the dots represent the 24 months of 2004-2005 and are delineated by color:
red = summer (May-September) and blue = winter (October-April).





1078 Figure 6. Normalized summary (a) target and (b) Taylor diagram illustrating model skill 1079 of MLD<sub> $\rho$ </sub> and MLD<sub>0</sub>, max d $\rho$ /dz, and max dO/dz at 13 Chesapeake Bay stations for 1080 2004-2005. The "x" represents the skill of a model that perfectly reproduces the 1081 observations. Since RMSD of stratification is only computed at stations where both the 1082 observations and model exhibit stratification, the Model Mean is not calculable for these 1083 variables.





Figure 7. Normalized target diagrams for Models A-H demonstrating the (a) temporal and (b) spatial skill in resolving the variability of the strength of density stratification (circles) and the depth of pycnocline initiation (diamonds). In (a) the individual dots represent the 13 stations along the main stem of the Chesapeake Bay. In (b) the dots represent the 24 months of 2004-2005 and are delineated by color: red = summer (May-September) and blue = winter (October-April).

1095



1097 Figure 8. Normalized summary (a) target and (b) Taylor diagram illustrating model skill

- 1098 of bottom temperature, salinity, chlorophyll, nitrate, and dissolved oxygen at 13
- 1099 Chesapeake Bay stations for 2004-2005. The "x" represents the skill of a model that
- 1100 perfectly reproduces the observations.
- 1101
- 1102
- 1103







1107 Figure 9. Time series of bottom dissolved concentrations for station CB4.1C. Red dots

1108 represent the 34 observations made during 2004-2005. Grey lines are the individual

1109 model simulations. The dark blue line represents the model mean while the cyan line

1110 represents the 95% confidence interval of the model simulations.





1121 Figure 10. Scatter plots comparing observations of (a) the strengths of stratification of the

1122 pycnocline and oxycline and (b) the oxygen- and density-defined mixed layer depths.

1123 Size of the circles is proportional to the number of observations. Observations are from

1124 1998-2006 at the 13 Chesapeake Bay stations shown in Figure 2.



1133 Figure 11. Time series of observations at Station CB4.1C from 2003 – 2006 for (a)

bottom dissolved oxygen, (b) dissolved oxygen at the MLD<sub>0</sub>, and (c) MLD<sub>0</sub>.



- 1147 Figure 12. Time series of observations of dissolved oxygen and MLD<sub>0</sub> contours at
- 1148 Station CB4.1C for 2004 and 2005.