

1           **Evaluating Current Statistical and Dynamical Forecasting Techniques for**  
2                                   **Seasonal Coastal Sea Level Prediction**

3  
4  
5 Xiaoyu Long<sup>a,b</sup>, Matthew Newman<sup>b</sup>, Sang-Ik Shin<sup>a,b</sup>, Magdalena Balmeseda<sup>c</sup>, John Callahan<sup>d</sup>,  
6 Gregory Dusek<sup>d</sup>, Liwei Jia<sup>e</sup>, Benjamin Kirtman<sup>f</sup>, John Krasting<sup>e</sup>, Cameron C. Lee<sup>g</sup>, Tong  
7 Lee<sup>h</sup>, William Sweet<sup>i</sup>, Ou Wang<sup>h</sup>, Yan Wang<sup>a,b</sup>, Matthew J. Widlansky<sup>j</sup>

8  
9                                   <sup>a</sup> CIRES, University of Colorado Boulder, Boulder, Colorado, USA

10                                  <sup>b</sup> NOAA Physical Sciences Laboratory, Boulder, Colorado, USA

11                                  <sup>c</sup> European Centre for Medium-Range Weather Forecasts, Reading, United Kingdom

12                                  <sup>d</sup> NOAA Center for Operational Oceanographic Products and Services, Silver Spring, Maryland, USA

13                                  <sup>e</sup> NOAA Geophysical Fluid Dynamics Laboratory, Princeton, New Jersey, USA

14                                  <sup>f</sup> Rosenstiel School of Marine, Atmospheric, and Earth Science, University of Miami, Miami, Florida, USA

15                                  <sup>g</sup> Kent State University, Department of Geography, ClimRISE Laboratory, Kent, Ohio, USA

16                                  <sup>h</sup> Jet Propulsion Laboratory, California Institute of Technology, Pasadena, California, USA

17                                  <sup>i</sup> NOAA National Ocean Service, Silver Spring, Maryland, USA

18                                  <sup>j</sup> Cooperative Institute for Marine and Atmospheric Research, School of Ocean and Earth Science and  
19   Technology, University of Hawai'i at Mānoa, Honolulu, Hawai'i, USA

20  
21 Corresponding author: Xiaoyu Long, [xiaoyu.long@noaa.gov](mailto:xiaoyu.long@noaa.gov)  
22

23

## ABSTRACT

24       The need for skillful seasonal prediction of coastal sea level anomalies (SLAs) has become  
25 increasingly evident as climate change has increased coastal flooding risks. Here, we evaluate  
26 nine current forecast systems by calculating deterministic and probabilistic skill from their  
27 retrospective forecasts (“hindcasts”) over 1995-2015, for lead times up to 6-9 months, at two  
28 United States tide gauge stations (Charleston, SC and San Diego, CA). Additionally, we assess  
29 local skill enhancement by two post-processing/downscaling techniques: an observationally-  
30 based multivariate linear regression and a hybrid dynamical approach convolving sea-level  
31 sensitivity to surface forcings with predicted surface forcing variations. We find that all these  
32 approaches face challenges stemming from whether modeled SLAs sufficiently represent  
33 observed local coastal SLA variations, because of both ocean model limitations and  
34 inadequacies in model initialization and ensemble spread. Some of these issues also complicate  
35 the ability of the post-processing techniques to improve probabilistic skill, even though they  
36 do somewhat improve deterministic skill. In general, deterministic hindcast skill is  
37 considerably higher for San Diego than Charleston, as expected from the stronger influence of  
38 ENSO. However, ensemble spread metrics such as forecast reliability and sharpness remain  
39 low for both locations, highlighting model deficiencies in representing uncertainty.  
40 Additionally, evaluating how well any technique predicts seasonal coastal sea level variations  
41 is complicated by the forced trend component and particularly how it is estimated. Moreover,  
42 model skill is matched by a stochastically-forced multivariate linear prediction model  
43 constructed from observations, suggesting that substantial improvement remains for predicting  
44 coastal seasonal SLAs, which could also include leveraging other predicted fields including  
45 sea level pressure and prevailing winds.

46

47

## SIGNIFICANCE STATEMENT

48

Coastal floodings have occurred more frequently in the last few decades, and it is anticipated that the number of such hazardous events will be increasing in the future.

49

50

Therefore, accurate and reliable forecasting of coastal water level is becoming increasingly

51

more important. This study thoroughly evaluated some current forecast techniques for sea

52

level along the U.S. coasts and found that those techniques are still not capable to produce

53

useable forecasting of anomalous sea level at U.S. east coast 3 months in advance, due to

54

model inadequacy. The current generation of forecasting models were not designed for

55

coastal sea level prediction, and we propose a few potential improvements that can

56

potentially advance our capability in coastal sea level and inundation forecasting in the near

57

future.

58

## 60 **1. Introduction**

61 Coastal flooding is a growing concern for the United States (U.S.) due to ongoing sea  
62 level rise (Church et al., 2013; Church & White, 2006; May et al., 2023), seasonal-to-decadal  
63 sea level variability, and land subsidence (Nicholls et al., 2021). These flooding events  
64 impact both ecosystems and public safety, with damages to the natural environment and built  
65 infrastructure including backed-up drainages, road closures, and saltwater intrusions.  
66 Flooding risks are projected to continue increasing in the coming decades (Taherkhani et al.,  
67 2020; Thompson et al., 2021; Vitousek et al., 2017), so there is an urgent need for accurate  
68 and reliable coastal flood prediction, ideally for months to seasons in advance. High coastal  
69 water levels are forced by a combination of drivers occurring at multiple spatial and temporal  
70 scales, including tides, waves, anomalous river runoff, storm surges, and sea level anomalies  
71 (SLAs) driven by atmospheric and oceanic processes on both weather and climate time  
72 scales. While the tide is usually the largest component of local water level variability, coastal  
73 flooding events typically occur when high tide coincides with other conditions favorable to  
74 anomalously high water levels. Though the astronomical (tide) contributions to sea levels are  
75 already well predicted, a high tide flooding prediction system should consider all these other  
76 factors as well (Hague et al., 2023).

77 Dusek et al. (2022) recently proposed a statistical approach to predict subseasonal to  
78 seasonal high tide flooding for U.S. tide gauges that are part of NOAA's National Water  
79 Level Observation Network (NWLON), which predicts daily probabilities of exceedance of a  
80 predefined flooding threshold for each location. These forecasts were made by combining  
81 tide predictions with a statistical representation of the "non-tidal residual" component of local  
82 sea levels, which includes both the long-term linear trend and the climatological distribution  
83 of hourly SLAs. They also showed that these forecasts were improved by including a simple  
84 damped-persistence model of monthly SLAs based on the observed autocorrelation function  
85 determined for each location. We might expect further improvement by using climate model  
86 seasonal SLA predictions that are more skillful than damped persistence.

87 Numerous studies have linked SLAs to patterns of seasonal climate variability (Han et al.,  
88 2017; Han et al., 2019b; Long et al., 2020; Roberts et al., 2016; Wang et al., 2023). This  
89 suggests that SLAs might be potentially predictable on seasonal time scales (Shin &  
90 Newman, 2021), and many studies have assessed seasonal SLA prediction skill of both

91 dynamical and statistical models (Chowdhury et al., 2007; Long et al., 2023; McIntosh et al.,  
92 2015; Miles et al., 2014; Widlansky et al., 2017). Long et al. (2021) constructed a 10-model  
93 ensemble forecast and assessed its skill of forecasting SLAs with lead times up to 12 months,  
94 finding that the dynamical models produce more skillful forecasts than a damped-persistence  
95 model at most open ocean locations. However, these models generally do not have higher  
96 coastal skill than does an observationally-based Linear Inverse Model (LIM), a multivariate  
97 empirical dynamical model that also allows for transient anomaly growth (Shin and Newman  
98 (2021). Additionally, Frederikse et al. (2022) employed a hybrid dynamical approach, where  
99 observed surface forcings and predicted surface forcings, from hindcasts generated by state-  
100 of-the-art seasonal forecast models, were projected onto SLA sensitivity at a specified  
101 location to global surface forcings computed by an ocean adjoint model. The resulting SLA  
102 hindcasts for the Charleston, SC location compared more favorably to observed tide gauge  
103 values there than the SLAs predicted by the same forecast model. Complicating all these skill  
104 assessments is that the pronounced externally-forced trend in sea level provides a substantial  
105 component of skill, at least as measured using commonly-used metrics (Wulff et al. 2022),  
106 that obscures the models' ability to predict seasonal climate variations (Long et al., 2021;  
107 Shin & Newman, 2021).

108 All the above studies of seasonal SLA prediction primarily focused on deterministic  
109 forecasts (e.g., ensemble means) and their skill assessment. However, warning end-users  
110 about high-tide flooding risks requires information about the likelihood of high-water  
111 (extreme) events (Dusek et al. 2022), which entails predicting tail probabilities. Therefore, it  
112 is also important to assess the probabilistic skill of coastal SLA prediction. For climate  
113 models, differences between multiple ensemble members (i.e., multiple forecast realizations)  
114 capture how initial uncertainty impacts the relative likelihood of future climate states.  
115 Probabilistic skill assessment then becomes a comparison, over the entire hindcast period,  
116 between the predicted probabilities of some extreme event and the actual chances of  
117 observing that event.

118 To improve our ability to forecast coastal flooding risk on seasonal and longer time  
119 scales, NOAA and NASA initiated the RISE project, a collaborative effort focused on  
120 developing and assessing novel dynamical and statistical forecast methods of SLAs along  
121 U.S. Coasts. This paper is an outgrowth of that project, which initially focused on a pilot  
122 study of monthly SLA forecast skill for sample tide gauge stations on the U.S. West and East  
123 Coasts (San Diego, CA and Charleston, SC). In this study, we evaluate monthly hindcasts of

124 sea level anomalies for those two tide gauge stations using deterministic and probabilistic  
125 metrics. We also discuss challenges involved in making coastal SLA forecasts, including how  
126 trends in the model outputs impact skill assessment and how to use models that may not be  
127 correctly initialized with observed sea levels.

128 The paper is organized as follows. Section 2 reviews issues involved in making coastal  
129 sea-level predictions from the output of various dynamical and statistical models. Section 3  
130 describes the forecast and observational datasets and the general skill metrics and methods  
131 used in this study. Section 4 presents the results of the deterministic and probabilistic skill  
132 assessment of seasonal SLAs, including a discussion of how this skill could be considerably  
133 impacted by both the presence of the externally forced trend and the inability of some  
134 hindcasts to represent it. Some remarks on how forecast models that are not initialized with  
135 satellite altimetry might be corrected follow in Section 5. Concluding remarks are made in  
136 Section 6, including recommendations for future advances in seasonal forecast systems to  
137 improve our prediction of coastal SLAs.

## 138 **2. Challenges for Coastal Sea Level Seasonal Forecasts**

139 As introduced above, previous studies have examined coastal SLA seasonal prediction  
140 skill. Yet it remains unclear how climate model output of monthly sea surface height  
141 anomalies should best be used to predict the risks of coastal flooding, and especially how  
142 these predictions should be verified against sea levels that are observed at tide gauges along  
143 the U.S. coastline. In most coupled climate models, the global ocean volume is conserved  
144 (i.e., they employ the Boussinesq approximation; Griffies & Greatbatch, 2012). As a result,  
145 these models cannot represent the global increase in sea level due to steric (thermal  
146 expansion) or barystatic (changes in water mass) processes, although they do allow for local  
147 height changes due to vertically integrated divergence/convergence, which is reflected in the  
148 model “sea surface height” (variable “zos” in the output from the Coupled Model  
149 Intercomparison Project (CMIP), as described in Griffies et al. (2016). This approach is  
150 sufficient for some purposes because changes in the global mean volume do not impact either  
151 ocean dynamics or coupling to the atmosphere. While the dynamical models considered here  
152 use a non-linear free surface, some older models use a “rigid lid approximation” with no  
153 variations in sea level at all; in this case, sea level must be calculated diagnostically from the  
154 model’s ocean bottom pressure and density profiles (Griffies & Adcroft, 2008). In all cases,  
155 however, local density variations due to temperature and salinity changes can impact sea

156 level locally, under the restriction that their global integral remains constant under the  
 157 Boussinesq approximation.

158 The local change of sea level  $\eta$  can be expressed as (Griffies & Greatbatch, 2012):

$$159 \quad \frac{\partial \eta}{\partial t} = \frac{Q_m}{\rho(\eta)} - \nabla \cdot \mathbf{U} - \int_{-H}^{\eta} \frac{1}{\rho} \frac{d\rho}{dt} dz, \quad (1)$$

160 where  $Q_m$  is water mass flux from the boundary (land ice melting, river runoff, etc.),  
 161  $\rho(\eta)$  is the water density at the surface,  $H$  is depth of the ocean, and  $\mathbf{U}$  is the vertically  
 162 integrated ocean current vector. The first term on the right-hand side of (1) is the contribution  
 163 to the local sea level from mass input to the ocean, which is important to the inter-annual  
 164 variability of global mean sea level (Hamlington et al., 2020). The second term is the  
 165 vertically integrated divergence, accounting for the local change of dynamic sea level. The  
 166 last term is the non-Boussinesq steric effect, which arises from density change following a  
 167 fluid parcel and vanishes in a Boussinesq fluid. The global mean sea level change due to this  
 168 non-Boussinesq steric effect can be corrected diagnostically as (Griffies & Greatbatch, 2012):

$$169 \quad \eta(s, t) = \eta^B(s, t) + \frac{V_0}{A} \ln \frac{\rho(0)}{\rho(t)} \quad (2)$$

170 where  $\eta(s, t)$  is the sea level at a given location  $s$  at time  $t$ ,  $\eta^B(s, t)$  is dynamic sea level  
 171 from the ocean model,  $V_0$  is the initial reference volume of the global ocean,  $A$  is the global  
 172 surface area of the ocean,  $\rho(0)$  is the initial global volume-averaged density, and  $\rho(t)$  is the  
 173 global volume-averaged density at time  $t$ . To verify sea level forecasts against global  
 174 altimetry observations and reanalyses, and ultimately predict coastal sea level, this quantity,  
 175  $\eta(s, t)$ , needs to be predicted, but unfortunately many seasonal forecast systems typically  
 176 output only  $\eta^B(s, t)$ .

177 Global reanalyses that assimilate both in situ and satellite observations are used to  
 178 initialize and verify seasonal forecasts produced by climate models. For sea surface height,  
 179 global observations are only available since 1993, from satellite altimetry. Some ocean  
 180 reanalysis systems also use altimetry to correct the global mean sea level change (Balmaseda  
 181 et al 2013). However, not all seasonal forecast climate models include altimetry in their  
 182 initialization. Long et al. (2021) noted that models that included assimilation of altimetry data  
 183 in their initialization tend to have a more realistic trend (due to both internal variability and  
 184 external forcing) in their hindcasts than models that did not, which complicated skill  
 185 comparison between models, especially in regions of strong sea surface height trends such as  
 186 the U.S. East Coast. Initialization that uses altimetry appears to improve seasonal SLA

187 prediction in many ocean regions, although less obviously so along the North American  
188 coastline (Widlansky et al., 2023) where the benefit of altimetry observations may not have  
189 been fully realized in the present-generation of assimilation systems (Feng et al. 2024).

190 Complicating matters further is that station-based tide gauges measure the water level  
191 relative to benchmarks on land (Gill & Schultz, 2001; Pugh & Woodworth, 2014). Also,  
192 some gauges are in bays or inlets, which can complicate their relationship to coastal sea level,  
193 and their measurements may include effects of freshwater flows from upstream (Piecuch et  
194 al., 2018). Also, since the land itself may move over time, tide gauge measurements can  
195 implicitly include a component due to vertical land motion (VLM) (Wöppelmann & Marcos,  
196 2016). While for seasonal forecasts VLM is so small that it can be neglected (and none of the  
197 regional or global models simulate VLM), it can become important over the entire multi-  
198 decade period common to most seasonal hindcast datasets, where it is not easily accounted  
199 for (Ray et al., 2023; Zervas et al., 2013).

200 Seasonal climate model forecasts are often “mean bias-corrected”, a post-processing step  
201 in which potentially erroneous model climatological mean states are replaced with the  
202 observed climatological mean state, which typically depends upon both the seasonal cycle  
203 and the forecast lead time (Stockdale et al 1993). That is, forecasts are verified by  
204 comparison of observed anomalies (relative to the observed mean state) to predicted  
205 anomalies (relative to the model mean state at that lead time). Typically, mean states are  
206 defined over a few decades, long enough to reduce sampling effects but still reasonably short  
207 enough to be representative of the current climate state in the context of long-term (i.e.,  
208 centennial scale) climate change. Unfortunately, as climate warming has accelerated over the  
209 latter half of the 20th century, this latter assumption is not valid for SLAs at many locations,  
210 especially along the East and Gulf coasts that have experienced an accelerating trend in mean  
211 sea levels over the past few decades relative to the global mean (Hamlington et al., 2020).  
212 The presence of a trend, even over a relatively short climatological period, leads to two  
213 issues. First, as noted above, the forecast model may not be able to entirely simulate all the  
214 processes responsible for the sea level trend itself. Second, anomalies are typically defined  
215 relative to a fixed long-term mean, which means that the trend component is included as part  
216 of the anomaly and, therefore, has a pronounced impact on the estimation of seasonal skill  
217 (e.g., Wulff et al., 2022). This is illustrated in Fig. 1 by considering a simple case where  
218 hindcasts of seasonal variations are so unskillful that they are entirely uncorrelated with the  
219 observed time series, but where both hindcasts and observations are also superposed about a



220 common linear trend. Then, the resulting two time series would be well correlated (Fig. 1a),  
221 making the seasonal forecast system appear skillful. Conversely, an incorrect trend in the  
222 hindcasts relative to observations could reduce skill even where the seasonal variations of the  
223 hindcasts and observations were otherwise well correlated (Fig. 1b). In cases such as  
224 illustrated in Fig. 1, we might simply remove or correct the linear trend. For example,  
225 Balmaseda et al (2024) show that a simple linear trend correction adds skill to seasonal SLA  
226 forecasts. More generally, however, evaluating the impact the trend on local hindcast skill  
227 becomes problematic when the externally-forced trend is nonlinear, particularly for relatively  
228 short records when it is unclear how to separate the trend from natural variability (e.g.,  
229 Solomon et al. 2011).

### 230 **3. Data and Methods**

231 Here, we discuss how we assess hindcast skill from various seasonal forecast systems,  
232 encompassing purely dynamical models, purely statistical models, and hybrid techniques.  
233 Since previous studies of (deterministic) sea level forecast skill all used different hindcast  
234 periods, we assess skill of all these techniques for 1995-2015, which is the common period  
235 for hindcast availability from all the forecast techniques.

#### 236 **3.1 Tide gauge verification data**

237 The verification data are based on monthly mean sea level data at the San Diego and  
238 Charleston NOAA NWLON tide gauges from 1995 to 2016, obtained from the Permanent  
239 Surface for Mean Sea Level (PSMSL; Holgate et al., 2013). SLAs are defined by removing  
240 the 21-year monthly mean climatology from each tide gauge time series. We limited our skill  
241 assessment to these two stations since the adjoint model of the Estimating Circulation and  
242 Climate of the Ocean (ECCO) system had only developed hindcasts there.

243 While we focus on SLA prediction at the tide gauges, we also discuss results from three  
244 other observationally-based gridded datasets: the SSALTO/DUACS multimission satellite  
245 altimetry dataset (Hauser et al., 2021; also known as AVISO in the literature) with a  $1/4^\circ$   
246 spatial resolution, and three ocean reanalyses, ORAS5 (Zuo et al., 2019) with a  $1/4^\circ$  spatial  
247 resolution, ECCO (Forget et al., 2015) with  $1/4^\circ$  spatial resolution, and GLORYS12 (Jean-  
248 Michel et al., 2021) with a  $1/12^\circ$  spatial resolution. As discussed in Section 2, the tide gauges  
249 measure quantities beyond what may be captured by global oceanic datasets, which is  
250 illustrated by comparing the gauge time series with the nearest grid value from the gridded  
251 datasets, shown in Fig. 2. Also shown is the correlation of the monthly gauge-located SLAs

252 of each of the gridded datasets and the tide gauge time series. All the gridded products  
253 capture the San Diego tide gauge reasonably well, but ORAS5 and ECCO capture only about  
254 half the monthly SLA variance observed at the Charleston tide gauge, which may be partly  
255 due to the relatively low weight that coastal data is given in their data assimilation systems  
256 (Feng et al. 2024). Finally, we also show linear and quadratic trend lines to each tide gauge  
257 record, determined by fitting a line or quadratic curve, respectively, to the data by minimizing  
258 the mean squared error. Note that there appears to be some upward trend over the period of  
259 record, which seems to have accelerated for both gauges after about 2011 when global mean  
260 sea surface temperatures also began to increase more rapidly (Garcia-Soto et al., 2021), so it  
261 is unclear how well either least-squares fit captures the externally-forced trend component.

### 262 3.2 Hindcast techniques and data

263 First, we consider hindcasts from three traditional assimilation-initialized seasonal  
264 forecast systems based on coupled dynamical models: CCSM4 (Community Climate System  
265 Model Version 4, Kirtman et al., 2014), SPEAR (Seamless System for Prediction and Earth  
266 System Research, Delworth et al., 2020; Lu et al., 2020), and ECMWF SEAS5 (Johnson et  
267 al., 2019). Apart from other modeling framework differences, including horizontal resolution  
268 (see Table 1), these forecast systems differ in how the ocean state is initialized and how the  
269 ocean model simulates global mean sea level evolution (although note that all the ocean  
270 models have a free surface):

271 (1) CCSM4 has the Parallel Ocean Program version 2 (POP2) model as its ocean  
272 component and is initialized with the Climate Forecast System Reanalysis (CFSR). Though  
273 CFSR captures the realistic variation of ocean heat content and hence the variation of SLA  
274 (Xue et al., 2011), the POP2 model by construction requires the global mean sea level to  
275 remain constant. Consequently, CCSM4 has no trend in its global mean sea level. Also,  
276 CCSM4 does not account for global mean variations in freshwater fluxes.

277 (2) SPEAR uses its ocean data assimilation to initialize its ocean model, the Modular  
278 Ocean Model Version 6 (MOM6). This data assimilation incorporates observed temperature  
279 and salinity profiles from ARGO based on the Ocean Tendency Adjustment (OTA) outlined  
280 in Liu et al. (2020). SPEAR does not explicitly simulate the global mean steric sea level  
281 evolution from internal changes in heat and salt due to the previously discussed limitations of  
282 a Boussinesq model. Still, unlike CCSM4, it does consider the imbalance of the hydrological

283 cycle of the climate system, which is accounted for within the topmost layer of the ocean  
284 model (Cazenave et al., 2012).

285 (3) SEAS5 also uses a Boussinesq ocean model, the Nucleus for European Modelling of  
286 the Ocean (NEMO) model. Unlike SPEAR, however, the ocean initial conditions of SEAS5  
287 include information from the altimeter observations via assimilation, including the global  
288 steric change (Zuo et al., 2019). Therefore, the SEAS5 forecasts of sea level will inherit the  
289 information from the sea level trend in their initial conditions and have a more realistic sea  
290 level trend both regionally and globally.

291 The impact of some of these configuration differences in each forecast system is seen  
292 when comparing observationally-based monthly anomalies of global-mean sea level,  
293 determined from the ORAS5, GLORYS12, and AVISO datasets (Section 3.1), to globally  
294 averaged lead-1 month SLAs from the CCSM4, SPEAR, and SEAS5 hindcasts (Fig. 3).  
295 [Note that here we define the lead-1 month as the first month after initialization, so that it is a  
296 combined representation of the initial climate state and the short-term model evolution from  
297 it.] Figure 3a shows that while there are some differences between the observationally-based  
298 time series, all three capture the long-term trend in the global mean, as do the SEAS5 lead-1  
299 hindcasts. The observed evolution is not captured by the lead-1 hindcast anomalies output  
300 from either the SPEAR or the CCSM4 (Fig. 3b), which, as noted above, do not include the  
301 global-mean steric or barystatic components. For the SPEAR, following the approach  
302 discussed in Section 2, the global-mean steric component was computed from one ensemble  
303 member of the lead-1 hindcasts and added to the global mean lead-1 ensemble-mean hindcast  
304 (Fig. 3b) to produce the “SPEAR+steric” curve in Fig. 3a, yielding a closer match to  
305 observations. Still, there remains a discrepancy, likely due to the lack of information about  
306 changes in initial oceanic volume that assimilation of satellite altimetry could provide.  
307 Finally, comparing the linearly detrended global-mean sea level anomalies for the three lead-  
308 1 hindcast datasets to that determined from ORAS5 (Fig. 3c) shows the CCSM4 and SPEAR  
309 initializations may not entirely capture interannual variations in the global mean sea level,  
310 even apart from the trend.

311 We also created a downscaled version from each of the dynamical forecast ensembles,  
312 using the technique demonstrated in Long et al. (2023): A seasonally invariant deterministic  
313 downscaling operator was constructed by multivariate linear regression of the high-resolution  
314 ( $1/12^\circ$ ) GLORYS12 ocean reanalysis data against its coarse-grained ( $1^\circ$ ) counterpart. Then,  
315 the downscaling operator is applied to each ensemble member of each dynamical hindcast,

316 generating an ensemble of high-resolution coastal sea level hindcasts that we refer to as  
317 DownscalingCCSM4, DownscalingSPEAR, and DownscalingSEAS5, respectively. By  
318 downscaling each ensemble member rather than the overall ensemble-mean (as was done in  
319 Long et al. (2023)), we generate a downscaled hindcast ensemble whose spread is based upon  
320 the original model ensemble. Note that the downscaling operator is in a reduced Empirical  
321 Orthogonal Function (EOF) space, so that not all the variance of the original hindcasts is  
322 retained in the downscaled hindcasts. Multi-model ensemble means were constructed using  
323 either the hindcast ensembles from the three GCMs or their corresponding downscaled  
324 hindcast ensembles.

325 Frederikse et al. (2022) developed a hybrid dynamical approach for seasonal SLA  
326 prediction. They first computed the sensitivities of the coastal sea level at a specific location  
327 to different global atmospheric surface forcings, using the ECCO adjoint model. Then, SLA  
328 prediction is made by convolving these sensitivities to observed and predicted atmospheric  
329 surface forcings, made up of observed (ECCO) forcings up to 12 months prior to  
330 initialization time followed by predicted atmospheric forcings up to 12 months after  
331 initialization time. Note that the precise length of applied forcings depends upon forecast lead  
332 time (e.g., for a 5-month lead, the forcing consists of 12 months of observed forcing followed  
333 by 5 months of predicted forcing). In Frederikse et al. (2022), the predicted atmospheric  
334 forcing fields are from a 10-member CCSM4 model. In the present study, we also use a 15-  
335 member SPEAR model in addition to the 10-member CCSM4 model for the predicted surface  
336 forcings. The resulting hindcast ensembles are named ECCO\_CCSM4 and ECCO\_SPEAR,  
337 respectively. Hereafter, this approach is referred to as the ECCO adjoint approach.

338 Note that for both the CCSM4 and SPEAR, we assess the skill of the original dynamical  
339 hindcasts, the downscaled version of those hindcasts, and the ECCO adjoint model forced by  
340 those hindcasts (albeit using predicted atmospheric surface forcing variables rather than  
341 predicted sea surface heights). This yields an ideal suite of forecasts to compare each method  
342 because they are all derived from the same dynamical forecast system (CCSM4 or SPEAR).

343 Finally, we also included hindcasts from a LIM, trained using near-global gridded fields  
344 of SST from HadISST (Kennedy et al., 2019) and SLA from ORAS4 (Balmaseda et al.,  
345 2013) from 1961 to 2015 (Shin & Newman, 2021). The LIM's deterministic forecast is  
346 represented by its ensemble mean, and its fixed but lead-dependent expected error statistics  
347 are used to estimate the uncertainty (i.e., ensemble spread) of its forecasts (equation (8) in  
348 Penland & Sardeshmukh, 1995). The LIM hindcasts were ten-fold cross-validated for the

349 entire 1961-2015 period, but for this paper, we assess skill only for those hindcasts initialized  
350 in the common 1995-2015 period.

351 Similar to earlier studies (Frederikse et al., 2022; Long et al., 2021), we use a univariate  
352 AR1 model, or damped persistence (van den Dool, 2006), as a minimum baseline of skill for  
353 all the evaluated forecast techniques. Note that LIM and damped persistence are similar in  
354 that both are determined from the lead-1 autocovariance of the data, but the LIM yields a  
355 multivariate matrix operator rather than a univariate scalar, so it also yields transient anomaly  
356 growth that leads to additional state-dependent predictability (Shin & Newman, 2021). Like  
357 the LIM, an AR1 model has fixed but lead-dependent expected error statistics, which can be  
358 used to estimate its prediction uncertainty. The damped-persistence coefficients were  
359 determined from each tide gauge record, using data only during the hindcast period, and  
360 cross-validated using a leave-one-out methodology.

361 Most dynamical models are initialized with a near-instantaneous or daily field; this is on  
362 the first day of the month for the three dynamical forecast systems assessed here. The first  
363 monthly mean forecast is then the mean of the first month of the forecast run, sometimes  
364 called the “Month 0.5” forecast, i.e., centered in the middle of the calendar month (e.g.,  
365 Kirtman et al., 2014). In contrast, empirical models may be initialized with observed (tide  
366 gauge/reanalysis) monthly mean anomalies centered on the previous month, so that the 1-  
367 month lead LIM/damped-persistence forecast and the dynamical model Month 0.5 forecast  
368 verify simultaneously. Following Newman and Sardeshmukh (2017), for clarity we renamed  
369 both these forecasts the “Month 1” forecast (i.e., the first month of the forecast period), and  
370 so on for increasing forecast leads (see also schematic in Ding et al., 2018).

371 For the three dynamical forecast systems (CCSM4, SPEAR, and SEAS5) that are  
372 initialized using full-field variables, a mean bias correction is first applied by removing the  
373 lead-time dependent climatology determined during 1995--2015 (Smith et al., 2013), as  
374 discussed in Section 2. The statistical downscaling is applied to these bias-corrected anomaly  
375 fields. LIM hindcast anomalies, initially defined relative to the 1960-2015 period, are  
376 adjusted to be relative to the 1995-2015 climatology, but otherwise are uncorrected. The  
377 ECCO\_CCSM4 and ECCO\_SPEAR each are forced with the mean bias-corrected dynamical  
378 model atmospheric forcing ensemble members from CCSM4 and SPEAR hindcasts,  
379 respectively, with an additional mean bias-correction applied to the resulting adjoint model  
380 hindcasts.

### 381 3.3 Prediction skill metrics

382 Deterministic skill is assessed using the anomaly correlation coefficient (ACC) between  
383 observations and ensemble-mean predictions, as a function of lead time, either computed over  
384 all calendar months or calculated separately for each verification calendar month. ACC  
385 measures how well a model can predict the phase and sign of observed anomalies (Wilks,  
386 2011). We also computed the root-mean-squared (RMS) skill score (RMSSS; e.g., Newman  
387 and Sardeshmukh (2017)), defined as  $\varepsilon \equiv 1 - \hat{\sigma}$ , where the standardized error  $\hat{\sigma} = \sigma/\sigma_{obs}$ ,  $\sigma$   
388 is the RMS forecast error between observations and ensemble-mean predictions, and  $\sigma_{obs}$  is  
389 the observed climatological RMS value. RMSSS is a measure of the average relative  
390 amplitude of the forecast error, defined so that a perfect forecast has RMSSS=1, a  
391 climatological forecast (i.e., a predicted anomaly of zero) has RMSSS=0, and a forecast  
392 poorer than climatology has a negative score.

393 Probabilistic skill is assessed using two different metrics. First, we determined reliability  
394 diagrams (Weisheimer & Palmer, 2014), where hindcasts are grouped into bins according to  
395 the predicted probability (horizontal axis), and then plotted against the frequency at which  
396 observed events occur (vertical axis). For a perfectly reliable forecast system, predicted  
397 probabilities should match observed probabilities, in which case the reliability curve lies  
398 along the diagonal: If an event is predicted as having an x% probability of occurring, then the  
399 event should occur x% of the time. Also included are “sharpness” diagrams, showing how  
400 often each forecast probability is issued, particularly distinguishing forecasts other than the  
401 climatological probability (Wilks, 2011). For example, for a three-category tercile forecast, a  
402 sharp forecast system should be able to issue forecast probabilities other than the  
403 climatological probability of 0.33. We also calculate the reliability value as a single metric  
404 for easy comparison across techniques (Toth et al., 2006).

405 Finally, ROC (Receiver Operating Characteristic; Kharin & Zwiers, 2003) curves were  
406 constructed by plotting the false alarm rate against the hit rate for different probability  
407 thresholds. In general, as we lower the probability thresholds, more ‘positive’ forecasts will  
408 be issued, and hence, both the hit rate and the false alarm rate will increase. A good forecast  
409 system has a high hit rate while minimizing false alarms, so an ideal ROC curve is away from  
410 the diagonal towards the upper left corner of the diagram (Mason & Graham, 1999), thereby  
411 maximizing the area under the ROC curve (ROC area). ROC skill score (ROCS), defined as  
412  $ROCS = 2(\text{ROC area} - 0.5)$ , measures this quantity. ROCS values can range from -1 to +1,  
413 where  $ROCS < 0$  indicates skill worse than climatology (i.e., random chance).

## 414 **4 Hindcast skill**

415 Before we show the skill evaluation, in Figure 4 we display an example of hindcast  
416 ensembles from each of the techniques for both San Diego and Charleston, initialized as of  
417 August 1, 1997 (the LIM is initialized using the monthly anomalies of July 1997) with lead  
418 times up to 12 months. In this case, August 1997 is Month 1, September 1997 is Month 2,  
419 and so on.

420 The anomalous sea level at San Diego was heavily influenced by that year's strong El  
421 Niño event and the associated coastally trapped Kelvin wave that propagated along the west  
422 coast of North America from the Tropics (Hamlington et al., 2015; Ryan & Noble, 2002).  
423 This led to an observed SLA maximum of nearly 18 cm in San Diego by November. While  
424 the models had some hint of this coastal Kelvin wave, it was too weak and delayed by a  
425 couple of months (Balmaseda et al., 2002), which resulted in a relatively flat SLA response,  
426 even for the individual ensemble members. This error was also apparent for hindcasts  
427 initialized earlier in June and July, and it was not until the October initialization that the  
428 models captured the timing of the November maximum (not shown, but see  
429 <https://www.psl.noaa.gov/forecasts/SeaLevel/#RISE> for all hindcasts and verifications used  
430 in this paper). Interestingly, the three models also appeared to predict a similar delayed and  
431 too-weak SLA response in San Diego for the 2009-10 and 2015-16 El Niño events (not  
432 shown). The ensemble spread of the downscaled and ECCO adjoint approaches was mostly  
433 reduced compared to the original models' spread. The observed November maximum was not  
434 included within any technique's ensemble spread, including the LIM's 2 standard-deviation  
435 ensemble spread (shading).

436 For Charleston, all nine techniques predicted a flat SLA response with increasing lead.  
437 Notably, the observed February 1998 SLA maximum of 20 cm was not contained within the  
438 ensemble spread of any forecast technique. There is a striking difference between the  
439 observed tide gauge value and the Month 1 forecast for SPEAR, DownscalingSPEAR,  
440 SEAS5, and DownscalingSEAS5. The ECCO adjoint approach appears to have improved the  
441 comparison between Month 1 and the corresponding observed tide gauge monthly mean  
442 anomaly. Finally, note ECCO\_SPEAR appears to reduce ensemble spread relative to SPEAR,  
443 but the opposite is true for ECCO\_CCSM4 compared to CCSM4.

### 444 4.1 Deterministic skill

445 The deterministic skill of all the techniques for the common hindcast period (Fig. 5) is  
446 considerably higher at San Diego than at Charleston, in agreement with previous studies  
447 (Long et al., 2021; Shin & Newman, 2021). For San Diego, SEAS5 and DownscalingSEAS5  
448 have the highest skill, exceeding the multi-model ensemble mean (MMM) skill largely  
449 because the CCSM4 skill is poor. The LIM has skill comparable to MMM, SEAS5, and  
450 SPEAR for shorter leads, but its skill degrades faster for longer leads. The downscaling  
451 technique improves of CCSM4 but not SPEAR skill, consistent with Long et al. (2023). The  
452 ECCO adjoint approach improves upon CCSM4 skill even more than the downscaling  
453 technique but worsens SPEAR skill. The multi-model mean of the downscaled hindcasts has  
454 a similar skill to that of the dynamical model hindcasts. Note that for Months 1-2, only  
455 SEAS5 and the LIM have a skill that exceeds damped persistence (gray background shading),  
456 but as lead time increases, more techniques show relatively greater skill.

457 At Charleston, in contrast, the LIM has the highest skill, with SEAS5 having the highest  
458 skill of the dynamical models, slightly exceeded by the multi-model mean at longer leads.  
459 The ECCO adjoint approach improved the skill of both models, especially the CCSM4, while  
460 the downscaling technique slightly improved the skill of SPEAR but not of CCSM4.

461 As discussed in Section 2, hindcast skill evaluation is complicated by the existence of a  
462 pronounced externally-forced sea level trend, which has pronounced regional variations (e.g.,  
463 cf. San Diego and Charleston in Fig. 2) that also are captured differently by the initializations  
464 of the different forecast techniques. For example, the higher Charleston skill for SEAS5 and  
465 the LIM might be due to their more accurate initialization of the externally-forced trend,  
466 although recall from Fig. 2 that the ORAS5 (and likewise the ORAS4, which the LIM is  
467 trained upon) still has some important differences from the tide gauge in Charleston. To  
468 evaluate the impact of this trend on hindcast skill, however, we would need to first  
469 distinguish it from natural internal climate variability, which is complicated by the short  
470 observational dataset (e.g., Deser et al., 2014; Frankignoul et al., 2017). In turn, whether the  
471 externally-forced trend is linear or nonlinear, which is unclear from Fig. 2, can impact  
472 estimates of internal variability.

473 How to precisely determine the externally forced trend is beyond the scope of this study  
474 (although see discussion in Shin & Newman, 2021), so here we instead test the sensitivity of  
475 our skill assessment upon the two different trend estimates shown in Fig. 2. Specifically, we  
476 first remove either the linear trend or the quadratic trend from both hindcasts and verification  
477 data, determined separately for each, and then recompute skill metrics of the resulting



478 detrended data. The results for linear (quadratic) detrending are shown in Figs. 5cd (Figs.  
479 5ef). Additionally, since a forced trend could contribute to persistence, we recomputed the  
480 damped-persistence models separately for each detrending.

481 All the techniques have considerably reduced skill for the detrended hindcasts (verified  
482 against detrended observations) at Charleston, consistent with earlier studies (Long et al.,  
483 2021; Shin & Newman, 2021), suggesting that much of the apparent skill on the U.S. East  
484 Coast is due to the pronounced trend there over the past three decades (Han et al., 2019a).  
485 That is, much of the Charleston skill is not from hindcasts that capture month-to-month  
486 variations of sea level but rather arises because both hindcasts and observational anomalies  
487 have a large trend component (relative to the constant climatological 1995-2015 mean) that  
488 artificially inflates estimates of the prediction skill of monthly variations (e.g., Fig. 1a). Note  
489 that while the qualitative impact of detrending on skill is similar for all the techniques, its  
490 greatest impact occurs for those techniques with relatively realistic initializations including  
491 realistic trends. However, determining the quantitative impact of the trend on skill seems  
492 sensitive to the assumed form of the trend. For example, linear detrending (Fig. 5d) causes a  
493 larger decrease in SEAS5 skill than does quadratic detrending (Fig. 5f). On the other hand,  
494 the skill of both SPEAR and the ECCO adjoint approaches decreases more for quadratic  
495 detrending. Still, for both detrending methods, the skill of many of the forecast techniques  
496 exceeds damped persistence even at longer forecast leads. Finally, note that while linear  
497 detrending has a much less pronounced impact on hindcast skill for San Diego, its impact is  
498 not negligible, especially for some techniques and with increasing leads.

499 Using RMSSS as the deterministic skill metric rather than ACC (Fig. 6) gives a similarly  
500 qualitative picture for San Diego skill, including the relative ordering of skill across the  
501 techniques and the impact of detrending. However, for Charleston, with the RMSSS metric,  
502 there are now fewer techniques whose skill exceeds damped persistence, although the relative  
503 position of the LIM is unchanged; after detrending, only a few techniques have skill that even  
504 matches damped persistence, with RMSSS that is only slightly positive. Additionally, it is  
505 apparent in Fig. 6 that when the dynamical models' ACC goes below about 0.4, RMSSS  
506 becomes negative, indicative of skill worse than a fixed prediction of a zero anomaly (i.e.,  
507 climatology). However, this is not the case for either the LIM or damped persistence; for  
508 example, for damped persistence, RMSSS approaches zero only as ACC approaches zero.  
509 Interestingly, for single-member forecast systems,  $ACC=0.4$  is equivalent to 100%  
510 standardized error (Livezey & Chen, 1983).

511 Hindcast skill for both stations depends upon the target month, especially for San Diego  
512 (Fig. 7). This result was found in some previous studies (Long et al., 2023; Shin & Newman,  
513 2021) as well, but here the seasonality of skill for all the forecast techniques is compared  
514 together for a common period. For San Diego, skill maximizes for forecasts verifying during  
515 winter, which could be associated with the El Niño-Southern Oscillation (ENSO) signal that  
516 typically also has a wintertime maximum along the West Coast (e.g., Shin & Newman,  
517 2021). Some techniques also show skill for up to 1-2 season leads when verifying during  
518 summer, consistent with an ENSO signal that is not predictable until after spring (the “spring  
519 predictability barrier”; e.g., Tippett & L’Heureux, 2020). Interestingly, the ECCO adjoint  
520 approach especially improves skill during summer, particularly for the CCSM4, whose skill  
521 is otherwise considerably worse than the other models. In contrast, changes in skill from  
522 downscaling had a much weaker seasonal dependence. Still, neither ECCO\_CCSM4 nor  
523 ECCO\_SPEAR summertime skill exceeds that of the LIM or SEAS5. SEAS5 shows  
524 significant skill up to Month 7 throughout the year. Linear and quadratic detrending have  
525 quantitatively similar effects to those in Fig. 5, mainly for spring and summer verifications  
526 when skill is already lower (not shown).

527 For Charleston, where skill is generally much lower than for San Diego, there is a less  
528 clear impact of seasonality upon skill (Fig. 8). Some of the techniques (DownscalingCCSM4,  
529 SPEAR, and DownscalingSPEAR) only have significant positive ACC during the early  
530 winter months. In contrast, SEAS5 and DownscalingSEAS5 have significant skill during  
531 spring and somewhat during fall, even at the longest leads available, with the LIM having a  
532 similar pattern with generally higher ACC values. Again, the ECCO adjoint approach boosts  
533 skill during the warm season, yielding significant skill for April to August verifications  
534 through Month 5 that exceeds (albeit not significantly) the skill from any other techniques.  
535 Interestingly, much of the ECCO adjoint skill increase for Charleston is retained even after  
536 detrending (Fig. 9). However, the degree of improvement is less when the trend line is  
537 quadratic (not shown) rather than linear.

## 538 4.2 Probabilistic skill

### 539 4.2.1 Reliability and sharpness

540 As was the case for deterministic skill, probabilistic metrics are better for San Diego (Fig.  
541 10) than for Charleston (Fig. 11). Reliability is generally better for predictions of upper than  
542 lower tercile events. The most reliable forecasts are made by the LIM and SEAS5, even for

543 the lower tercile events. For San Diego, the LIM is more reliable than the other techniques at  
544 this lead (Month 4), even as its deterministic skill is relatively poorer. For Charleston, except  
545 for the LIM, none of the hindcasts are particularly reliable, and even the LIM is more reliable  
546 for the upper than the lower tercile. The ECCO adjoint technique also has minimal impact on  
547 reliability, slightly improving CCSM4 but making SPEAR worse, notably for the lower  
548 tercile. The multi-model means do not improve overall reliability at either location (not  
549 shown) and slightly degrade it for the downscaled hindcasts. However, this may be due to the  
550 small number of models we used.

551 The inset sharpness diagrams show that while all the models can issue forecast  
552 probabilities other than the climatological value of 0.33, these hindcasts are dominated by  
553 very low forecast probabilities (i.e., they are generally clustered in the leftmost bins).  
554 Interestingly, SPEAR and SEAS5 tend to have more forecasts with a higher probability for  
555 Charleston than for San Diego (note their U-shape sharpness diagrams). Finally, note that the  
556 adjoint technique reduces the occurrence of higher forecast probabilities (e.g., 0.7 and 0.9  
557 bin) compared to the models upon which they are based, particularly for Charleston,  
558 consistent with the reduced ensemble spreads for the ECCO\_CCSM4 and ECCO\_SPEAR  
559 seen in Fig. 4.

560 After linearly detrending Charleston hindcasts and observations, hindcasts become much  
561 less reliable (Fig. 12). That is, much of the (limited) reliability seen in Fig. 11 represents  
562 probabilities from hindcasts either early in the hindcast period when both hindcast and  
563 observed anomalies include a trend component that is relatively large and negative (when  
564 defined as an anomaly relative to the long-term mean), or late in the period when that trend  
565 component is relatively large and positive. It is not entirely surprising that the reliability and  
566 sharpness of the techniques with more realistic initialization (SEAS5 and LIM) are most  
567 impacted by the detrending. As with the deterministic metrics, the impact of the trend on  
568 reliability is much less at San Diego. Additionally, reliability is less sensitive than ACC to  
569 removing a quadratic rather than a linear trend, at least for the three-category approach used  
570 here (not shown).

#### 571 4.2.2 ROC skill scores

572 The ROC curves for San Diego for Month 4 (Fig. 13) show that all the techniques have  
573 better performance in predicting upper tercile than lower tercile events. SEAS5, Downscaling  
574 SEAS5, and LIM have the best performance, with ROCS values between 0.78 to 0.82 for the

575 upper tercile and 0.65 to 0.77 for the lower tercile, followed in order by the SPEAR and then  
576 the CCSM4. The downscaling and ECCO hybrid methods improve probabilistic skill for  
577 CCSM4 but not so much for SPEAR. The multi-model mean improves upon CCSM4 and  
578 SPEAR (and their related hindcasts) for both upper and lower tercile hindcasts but does not  
579 improve upon SEAS5. Note also that many techniques have relatively low hit rates and false  
580 alarm rates even for the lowest classification threshold because the techniques do not issue  
581 enough ‘positive’ event predictions overall.

582 For Charleston at Month 4 (Fig. 14), ROCS values for upper and lower terciles are  
583 generally lower than for San Diego. The LIM, SEAS5, and DownscalingSEAS5, in that  
584 order, have the highest ROCS values. The CCSM4 and its derived models  
585 (DownscalingCCSM4 and ECCO\_CCSM4) all have ROCS near zero, indicating that the  
586 CCSM4 has no skill compared to a random classification model. The SPEAR and its derived  
587 models are only slightly better. The multi-model mean also again improves upon CCSM4 and  
588 SPEAR but not SEAS5.

589 These results are representative of other forecast lead times, which is demonstrated by the  
590 ROCS values for each of the techniques as a function of lead (Figs. 15 and 16). For San  
591 Diego (Fig. 15), SEAS5, DownscalingSEAS5, and LIM have the highest ROCS values  
592 through Month 7. The CCSM4 and its derived hindcasts have the lowest ROCS values, with  
593 SPEAR in between. Neither downscaling nor ECCO hybrid methods improve ROCS values  
594 of SEAS5 and SPEAR, although the ECCO hybrid does improve CCSM4. For Charleston  
595 (Fig. 16), the overall skill scores are lower than for San Diego, where again, the highest  
596 ROCS values are for SEAS5, DownscalingSEAS5, and LIM. Similarly, the downscaling  
597 technique does not improve the skill score, while the ECCO hybrid method slightly improves  
598 the skill over shorter lead times.

599 The removal of the linear trend also effectively reduces ROCS for SEAS5 and  
600 DownscalingSEAS5, but less so for the LIM (Fig. 16). Detrending improves the ROCS of  
601 CCSM4 and DownscalingCCSM4, mostly because the removal of the spurious trend in those  
602 model hindcasts improves the quality of their hindcasts, in the same way that the detrending  
603 degrades the hindcasts which have more realistic trends (Fig. 1b). In contrast to the  
604 deterministic skill, the impact of the trend upon probabilistic skill does not much depend  
605 upon whether the trend is estimated as quadratic or as linear (not shown).

## 606 **5 Correcting forecasts from models with inadequate sea level initialization**

607 As discussed in Section 2, using and evaluating dynamical model predictions of coastal  
608 SLAs can be challenging when the model output sea level variable does not entirely  
609 correspond to tide gauge measurements, particularly when the model is not initialized with  
610 sea level observations. Perhaps this could be alleviated by evaluating prediction skill after  
611 removing the global mean from both verifications and model hindcasts, addressing the  
612 absence of the steric contribution to global mean sea level in forecast models and/or their  
613 initializations. We considered this approach as a possible solution both to this problem and to  
614 the problem of evaluating detrended skill since a significant portion of the trend is related to  
615 global mean sea level rise (e.g., Fig. 3a), although for some, if not most, tide gauges the trend  
616 also has VLM contributions. Unfortunately, global mean sea level is also impacted by large  
617 scale internal variability (e.g., Fig. 3c), including potentially predictable climate variations  
618 such as ENSO (Cazenave et al., 2012; Wang et al., 2021). Hence, removing the global mean  
619 component was inadequate to comprehensively evaluate the skill of seasonal prediction of  
620 regional sea level anomalies. Likewise, removing basin-wide means was also unsuccessful.

621 Alternatively, we might assume that models whose sea level is incompletely initialized  
622 (e.g., without altimetry in the data assimilation) might be capable of predicting month-to-  
623 month sea-level changes so long as other ocean variable initializations (e.g., temperature and  
624 salinity) are not substantially impacted. Such an initialization error might be considered  
625 simply an offset, addressed by adjusting the hindcasts. Let the predicted monthly sea level  
626 state for initial time  $t$  and lead time  $j$  be  $\tilde{Z}(t, j)$ . Then the prediction increment, or “delta”, of  
627 sea level can be calculated from the model output as:  $\Delta\tilde{Z}(t, j) = \tilde{Z}(t, j) - \tilde{Z}(t, j - 1)$ , for all  
628 lead times. Finally, the adjusted forecast is determined by incrementing the observed initial  
629 monthly anomaly,  $Z(t, 0)$ , with each delta at different lead times (i.e.,  $Z(t, 0) + \Delta\tilde{Z}(t, j)$ ),  
630 yielding a “delta-corrected sea level prediction” whose month-to-month change is identical to  
631 the original model forecast but is “initialized” with observations. However, this leaves us still  
632 with the choice of the initial observed monthly SLA. The most recent observed monthly sea  
633 level at the tide gauge is the previous month, or  $Z(t, -1)$ . Using that as the initialization  
634 means that we still need  $\Delta\tilde{Z}(t, 0)$ , which we determine from the difference between the  
635 current Month 1 and previous Month 1 forecasts ( $\Delta\tilde{Z}(t, 0) = \tilde{Z}(t, 0) - \tilde{Z}(t - 1, 0)$ ); then,  
636  $\tilde{Z}(t, 0) = Z(t, -1) + \Delta\tilde{Z}(t, 0)$ , and the forecasts can be incremented from that point onwards.

637 There are two benefits of this correction. First, the inconsistent trend between the model  
638 and observations is no longer an issue because the realistic trend is built into the corrected sea  
639 level. Second, the imperfect initialization is less of a problem since starting from the previous

640 month's observed SLA will largely eliminate the difference between model initialization and  
641 observation.

642 The delta-corrected sea level prediction is particularly appealing for issuing real-time  
643 SLA predictions that are appropriate for specific tide gauges, especially when using models  
644 initialized without a trend component to predict SLAs for tide gauges with pronounced  
645 observed trends (which over a few decades could also include vertical land motion). Note that  
646 all delta-corrected tide gauge hindcasts include the observed trend after correction, allowing  
647 comparison with other techniques already initialized with observations. The delta-correction  
648 improves the ACC for the models without a correct trend, both the CCSM4 and the SPEAR,  
649 especially for San Diego (Fig. 17; cf. lines with circles to same-colored lines with crosses).  
650 Note that not only is the skill of the original model hindcasts improved, but the skill of the  
651 related downscaled and hybrid model hindcasts are as well. Interestingly, the results at  
652 Charleston are not as consistent, even though the sea level trend is larger at Charleston than  
653 San Diego. This may be due to the larger vertical land motion component at San Diego  
654 (Zervas et al., 2013), which the delta-correction could also capture.

655 Note that an error is introduced for the delta estimate in Month 1 since it uses two  
656 separate model hindcasts initialized at two different times. This error could then propagate  
657 through the delta-corrected hindcasts for all lead times. The delta-correction degrades the  
658 skill of both SEAS5 and the LIM (not shown), whose initializations include observed sea  
659 level information that better captures observed trends. Hence, our delta-correction method is  
660 only an interim remedy for models with inadequate initialization and/or that do not output the  
661 global mean steric component forecasts.

## 662 **6 Concluding remarks**

663 In this study stemming from the RISE project (a collaboration among scientists at NOAA,  
664 NASA/JPL, and several universities), we have considered some key issues in the prediction  
665 of coastal SLAs on seasonal time scales, a particularly challenging problem since seasonal  
666 forecast systems largely have not been designed with such predictions in mind. Using both  
667 deterministic and probabilistic metrics, we assessed the skill of hindcasts from various  
668 dynamical and statistical models/techniques — traditional assimilation-initialized seasonal  
669 forecast systems based on coupled dynamical models, an empirical regression-based  
670 approach (the LIM), and two statistical (linear regression) and dynamical (ECCO adjoint)  
671 post-processing techniques applied to output from the seasonal forecast models — against

672 monthly SLAs observed at two sample NOAA NWLON tide gauge stations in San Diego,  
673 CA and Charleston, SC. We found that the skill of some of the forecast systems cannot beat a  
674 simple “damped persistence” (univariate AR1) approach, especially for Charleston. Even  
675 fewer had deterministic or probabilistic skill greater than the LIM (multivariate AR1)  
676 approach, which suggests that for future studies, the LIM could serve as a more rigorous  
677 benchmark than damped persistence for coastal SLA seasonal forecast skill, both  
678 deterministic and probabilistic.

679 Consistent with previous studies (Long et al. 2021; Shin and Newman 2021), SLA  
680 seasonal prediction skill was considerably better for San Diego than Charleston. There are a  
681 few possible reasons for poorer Charleston skill. There may simply be lower inherent SLA  
682 predictability in the Charleston region. For example, past studies have shown that while  
683 ENSO drives a strong and potentially predictable signal in Pacific SLA along the U.S. West  
684 coast (Amaya et al. 2022), predictable SLAs along the U.S. East Coast appear associated with  
685 Gulf Stream modulation that may have a smaller impact on the predictable monthly signal,  
686 compared to unpredictable noise processes (Shin and Newman 2021). This difference is  
687 likely exacerbated by large-scale climate model errors in the position of the Gulf Stream, due  
688 in part to model grids that are too coarse (e.g., Bryan et al., 2007), although there are also  
689 systematic errors in ENSO prediction as well (e.g., Beverley et al. 2023). Current climate  
690 model resolution may also be insufficient to entirely capture climate-related signals  
691 propagating along the coast, including Kelvin waves driven by ENSO (e.g., Amaya et al.  
692 2022) and other coastally-trapped waves (e.g., Brunner et al. 2019 and references therein;  
693 Hughes et al. 2019). Finally, inadequate initialization is also likely a contributing factor in  
694 poorly performing forecasts, especially around the Gulf Stream region (Widlansky et al.  
695 2023). For example, the CCSM4 and SPEAR Month 1 hindcasts, which were not initialized  
696 using the altimetry observations, failed to correctly represent either the (relatively large) trend  
697 along the US East Coast or the interannual component of global mean SLA (Fig. 3). Reliance  
698 of the LIM and SEAS5 on ORAS4 and ORAS5, respectively, means that their Charleston  
699 initializations are also somewhat deficient relative to San Diego (Fig. 2); since these  
700 reanalyses often give little-to-no weighting to satellite altimetry data near the coast (e.g.,  
701 Balmaseda et al. 2013; Feng et al. 2023), this difference may be related to the fact that while  
702 the San Diego tide gauge SLA is correlated with a large scale North Pacific SLA pattern,  
703 Charleston is primarily correlated with SLA along the South Atlantic coast (Long et al.  
704 2023). The resolution issues extend to the verification process, where coarse-grid model

705 hindcasts (output at 1° grid resolution) are compared against point observations, rather than  
706 against similarly gridded observational datasets. Interestingly, the ECCO adjoint approach,  
707 which is “initialized” using a 12-month dynamical spin-up forced by surface observations,  
708 has much better Charleston skill compared to the climate models even for Month 1 hindcasts,  
709 with significant skill (even after detrending) up to about Month 6 during late spring and late  
710 fall. Further diagnosis of these issues, analyzing how inadequate initialization and model  
711 error interact with each other so that forecast systems may not take full advantage of sources  
712 of predictability, will also need to consider other US tide gauges as well.

713 A common problem for all the hindcasts is that they do not appear to generate enough  
714 categorical “hits”, even for categories with the lowest probability thresholds. This issue is  
715 evident in individual hindcast ensembles (e.g., Fig. 4), as the ensemble spread is often  
716 insufficient relative to observed variability so the model forecasts are not entirely reliable.  
717 Note that while we do not expect climate models to predict observed SLA evolution with  
718 deterministic certainty, we do expect them to be able to produce ensemble members that can  
719 encompass what is observed to occur, so this over-confidence of the model hindcasts is likely  
720 also reflective of model error. While both post-processing approaches (downscaling and  
721 ECCO adjoint) yield some improvement for deterministic skill relative to the original  
722 hindcasts upon which they were based (albeit not uniformly across the models), neither  
723 technique appears to improve probabilistic skill and, in some cases, may reduce skill by  
724 collapsing the ensemble spread. This suggests the importance of developing post-processing  
725 and downscaling methods, whether dynamical or machine learning-based, whose ensembles  
726 can capture variability more realistically on the local scales of interest, even if much of that  
727 variability is unpredictable.

728 In this study, we detrended observations and hindcasts using either a linear or quadratic  
729 fit to explore the externally-forced trend’s potential impact on hindcast skill, including  
730 whether different trend estimates could yield different impacts on skill. For Charleston, we  
731 found that the pronounced sea level trend increases apparent seasonal prediction skill,  
732 especially for hindcasts that are realistically initialized. Much of the hindcast reliability at  
733 Charleston also appears to correspond to a trend component. While a trend-related impact on  
734 skill also exists for San Diego, it is much less pronounced and complicated by vertical land  
735 motion over the length of the hindcast period. For both stations, detrended hindcasts  
736 generally remain more skillful than the corresponding damped persistence benchmarks, since  
737 the trend also increases apparent persistence. Note that, for the most part, we determined only



738 the qualitative impact of the trend, since these impacts changed when the trend was assumed  
739 to be linear or quadratic. A quantitative analysis would require estimating and removing the  
740 evolving externally-forced trend, which is complicated by the presence of internal climate  
741 variability, vertical land motion, and the apparent trend acceleration since 2011. Of course, if  
742 the assumed trend differs from the actual trend, then incorrect detrending could remove some  
743 potentially predictable component of natural seasonal-to-interannual variability. This would  
744 impact our skill estimates, including metrics of the model's ability to capture observed  
745 marginal and conditional probability distributions (e.g., Xu et al. 2022). We also found that  
746 removal of the trend can sometimes improve skill for hindcasts with an erroneous trend  
747 component, due to inadequate initialization of a sea surface height forecast variable that does  
748 not represent the total sea level. Our attempt to correct this issue, by using forecast output to  
749 predict month-to-month SLA changes rather than the monthly SLA values themselves, still  
750 suffers from an inability to cleanly separate the initial observed state into its trend and  
751 seasonal anomaly components, and is thus, at best, a temporary, ad hoc fix. In essence,  
752 hindcast skill assessment of U.S. coastal seasonal SLA prediction is also an externally-forced  
753 trend detection problem.

754 This paper represents a multipronged assessment of the skill for seasonal prediction of  
755 regional sea level anomalies, involving dynamical, statistical, and hybrid methods. We have  
756 tried to provide important information to the climate prediction community about the relative  
757 strengths and limitations of various approaches, highlighting the challenges of sea level  
758 prediction at sample U.S. coastal stations, and stressing important issues to consider when  
759 assessing and comparing sea level prediction skill. Our primary conclusion is that, for the  
760 most part, the current seasonal forecasting systems may not yet be fit for the purpose of  
761 making coastal sea level predictions in the regions considered here. It is apparent that making  
762 useful predictions of coastal SLA is a hard test for seasonal forecast systems, which helps  
763 identify needs for additional improvement in both climate models and their initialization. To  
764 make progress, we propose studies aimed at the following:

- 765 • Evaluating the extent to which higher model resolution, which could reduce large  
766 scale model errors such as Gulf Stream position and strength, will improve  
767 forecasts that also depend upon complex coastal geography and bathymetry.
- 768 • Developing non-Boussinesq ocean models that include the global ocean volume  
769 changes, both barystatic and thermosteric, which are important to local sea level

770 prediction, or alternatively developing models whose output include global  
771 changes in steric volume as an additional diagnostic.

772 • Understanding how best to initialize climate models (either Boussinesq or non-  
773 Boussinesq) so that their coastal forecasts may be best used, either directly or  
774 through post-processing and downscaling.

775 • Investigating methods to improve the reliability of climate model forecast  
776 ensemble spread (e.g., stochastic parameterization; Sardeshmukh et al. 2023).

777 • Diagnosing drivers of model error in hindcast ensembles, especially the rapid  
778 initial development of error (including the initialization drift of the mean dynamic  
779 topography of the ocean) that degrades SLA forecast skill.

780 Additionally, the issues discussed concerning how to post-process/downscale seasonal  
781 climate forecasts lead us to suggest studies focused on:

782 • Evaluating how climate model forecast ensembles may be used for driving  
783 smaller-scale, limited domain ocean models, with better coastal processes  
784 including tide, wave, and ocean dynamic effects, that are all necessary for  
785 providing actionable coastal information.

786 • Constructing forecast ensembles that can capture prediction uncertainty both at the  
787 large climate scales and at the smaller atmospheric and oceanic scales relevant to  
788 the coastal regions.

789 • Evaluating new empirical and machine learning approaches, designed both to  
790 post-process model forecasts and to make coastal SLA seasonal predictions  
791 outright, as alternative solutions while the suggestions above are evaluated.

792 *Acknowledgments.*

793 XL, MN, and SS acknowledge the support from NOAA cooperative agreements  
794 NA17OAR4320101 and NA22OAR4320151, and U.S. DoC/ NOAA/Bipartisan  
795 Infrastructure Law (BIL). MJW was supported by the NOAA Climate Program Office's  
796 Modeling, Analysis, Predictions, and Projections (MAPP) program through grant  
797 NA22OAR4310138. Part of this research was carried out at the Jet Propulsion Laboratory,  
798 California Institute of Technology, under a contract with the National Aeronautics and Space  
799 Administration (80NM0018D0004).

800

801

802 *Data Availability Statement.*

803 All hindcast and observational datasets discussed in this paper are available for viewing  
804 and downloading at <https://www.psl.noaa.gov/forecasts/SeaLevel/#RISE>.

805

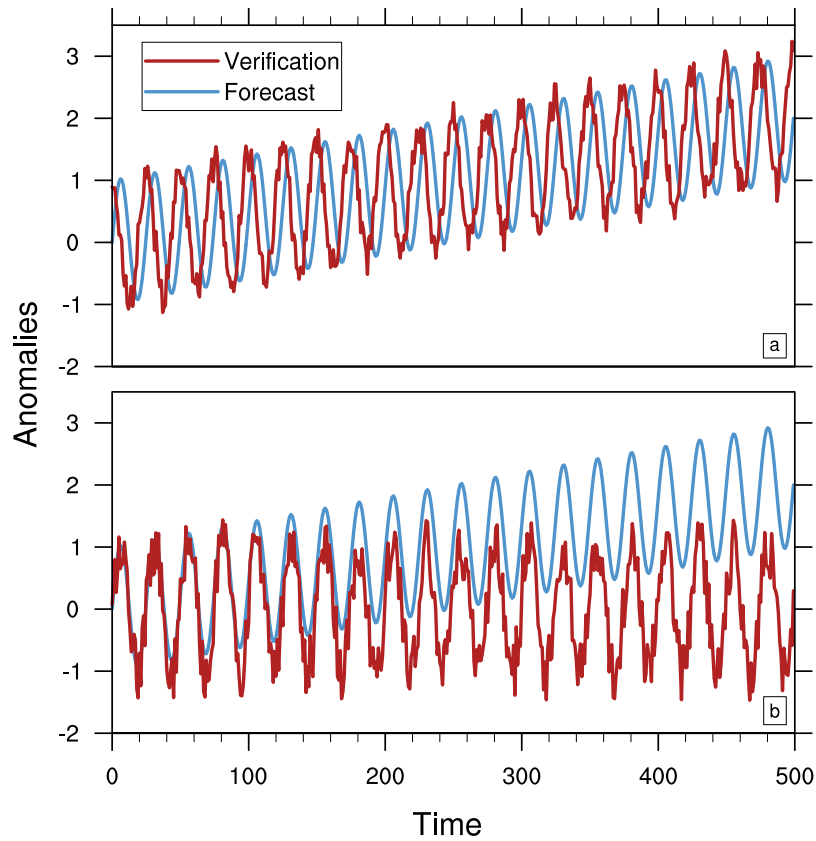
806 **Table 1** Key characteristics of the ocean model component in the three dynamical  
807 forecast systems, including the name of the forecast system, ensemble size of the seasonal  
808 hindcasts, lead time (months), the name of the ocean model component, nominal horizontal  
809 resolution, whether the system includes a global mean sea surface height component in its  
810 output, whether the system assimilates the altimetry observed sea surface height into its initial  
811 conditions, and the reference related to each of the seasonal forecast systems.

812

Name	Ensemble Size	Lead Time	Ocean Model	Grid Resolution (deg)	Global Mean Sea Level	Altimetry-initialized?	References
CCSM4	10	11	POP2	1	No	No	Kirtman et al. (2014)
SPEAR	15	11	MOM6	0.5	Partly	No	Delworth et al. (2020)
SEAS5	25	6	NEMO	0.25	Yes	Yes	Johnson et al. (2019)

813

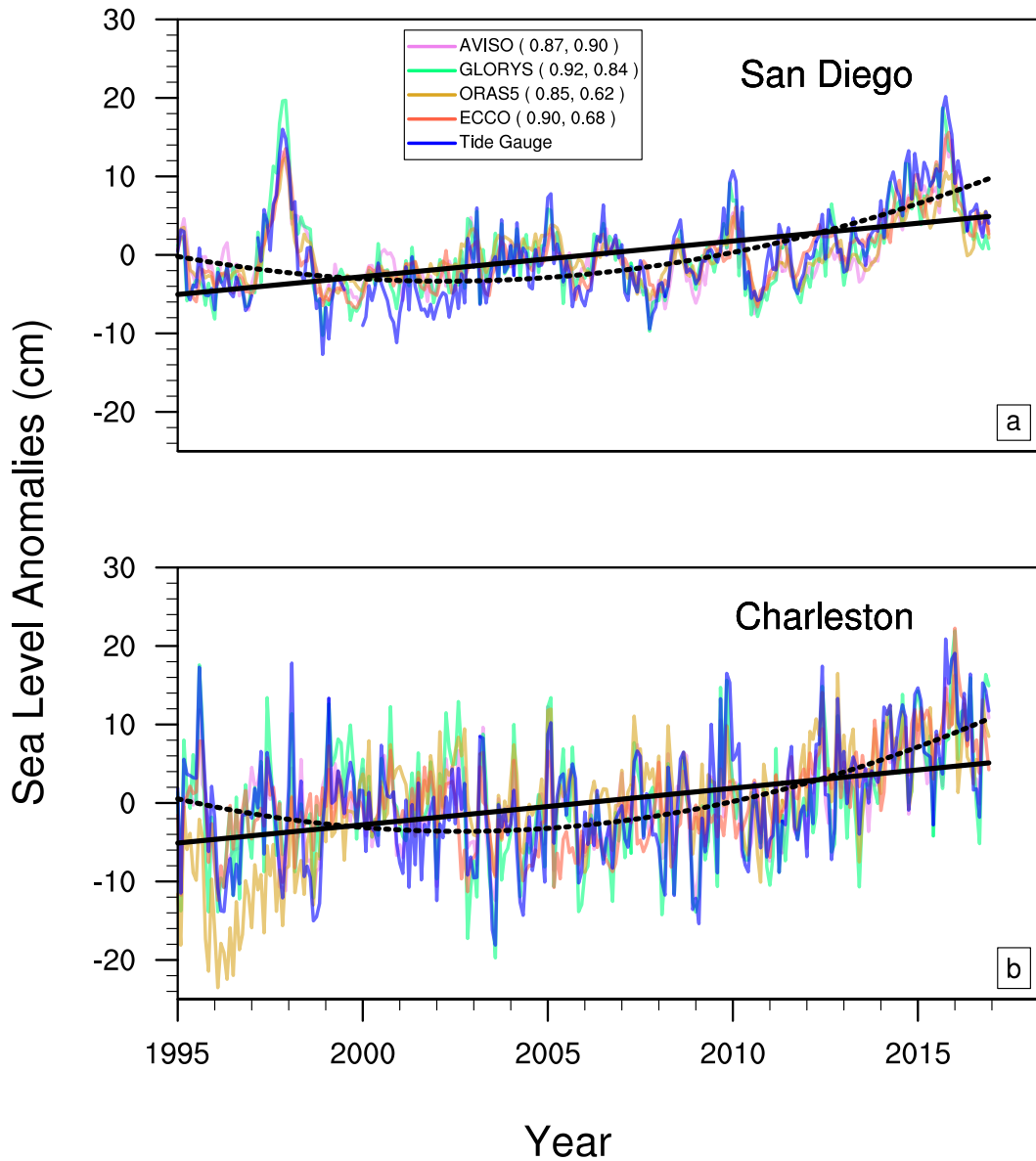
814



815

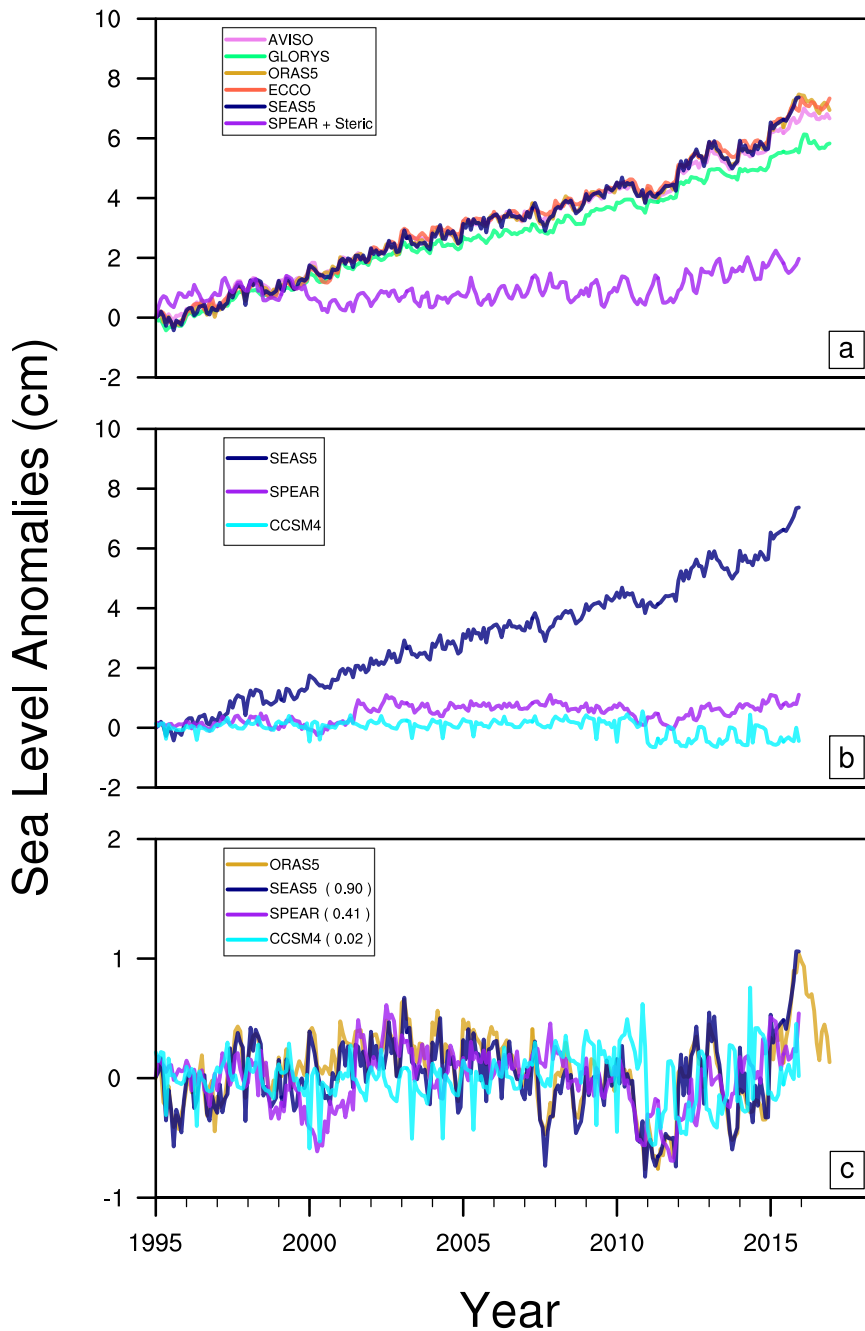
816 Figure 1. (a) An example of how superposing two uncorrelated time series onto a linear trend  
 817 generates the new well-correlated time series. The original uncorrelated time series has a  
 818 standard deviation of 0.9, and the superposed linear trend has a standard deviation of 0.56.  
 819 The two new time series correlate 0.39. (b) Superposing different trends can reduce the  
 820 correlation of two otherwise correlated time series from 0.90 to 0.67.

821



823

824 Figure 2. Monthly mean SLA from two tide gauge stations, (top) San Diego and (bottom)  
 825 Charleston, compared with SLA values of the nearest grid point from the observationally-  
 826 based reanalysis datasets, AVISO, GLORYS12, ECCO, and ORAS5. The correlation  
 827 coefficient between each reanalysis and the corresponding tide gauge time series is given as  
 828 numbers in parenthesis (San Diego, Charleston). Linear and quadratic least square fits to each  
 829 tide gauge time series are also shown.



830

831

832

833

834

835

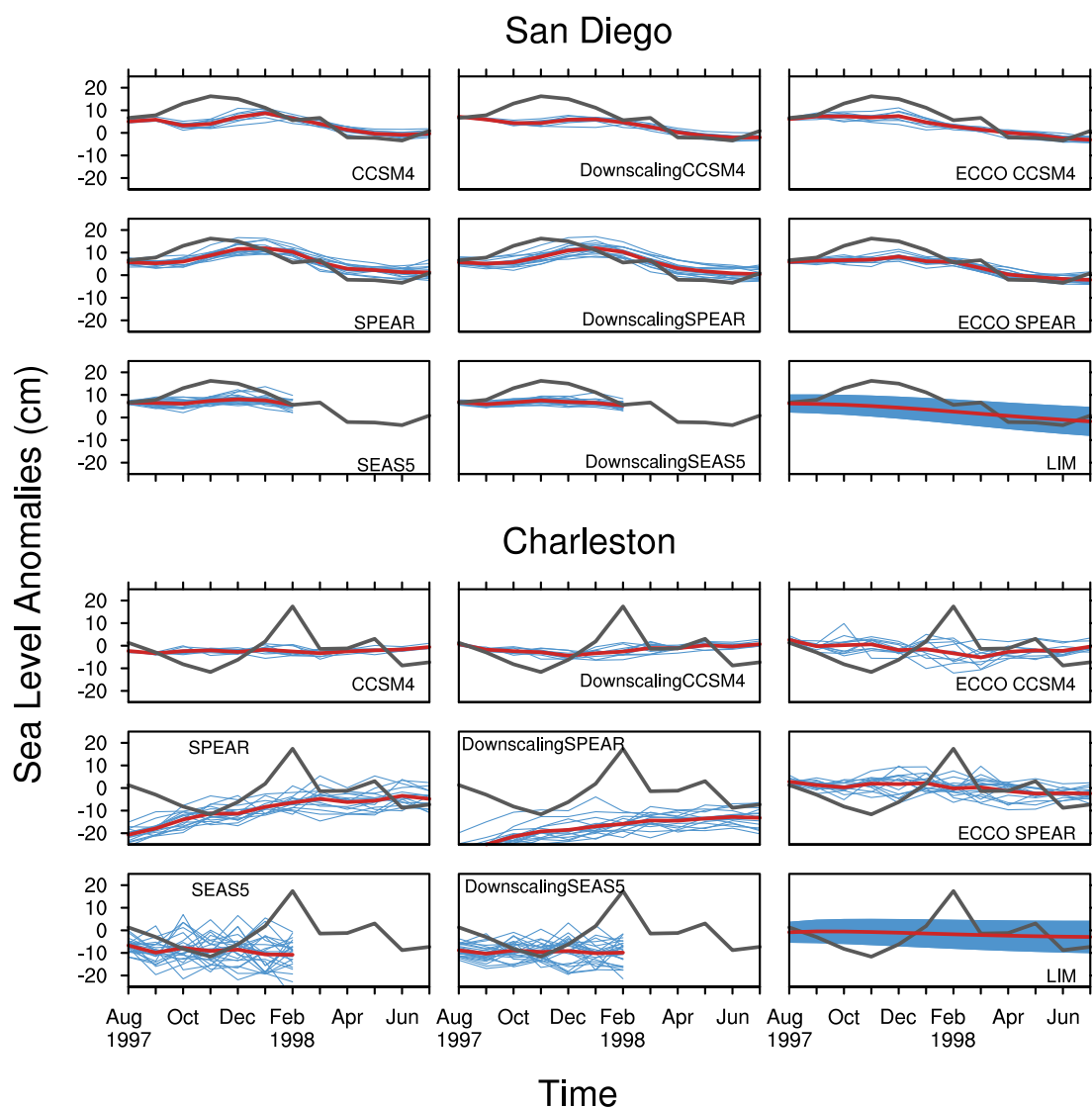
836

837

838

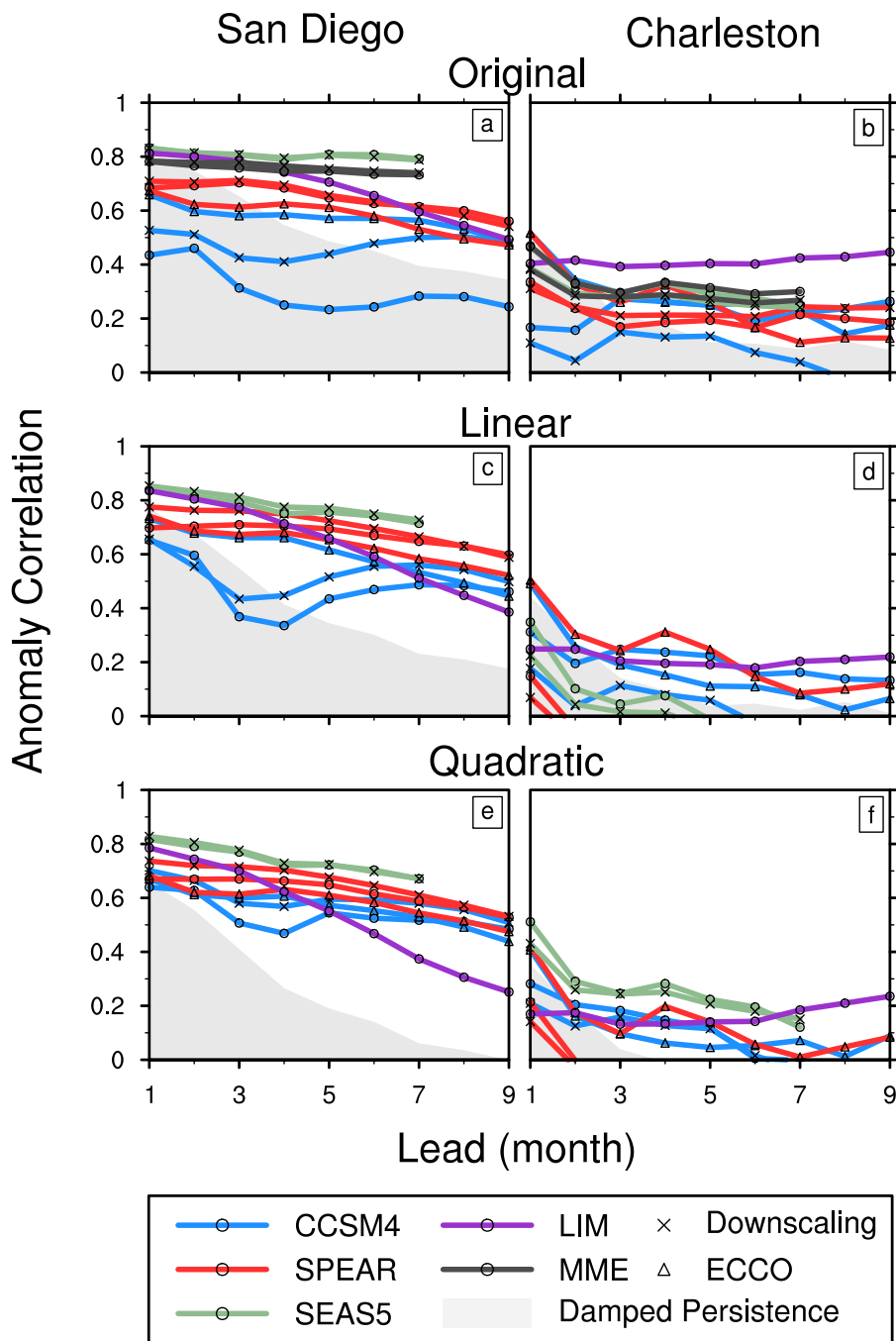
Figure 3. Global monthly mean SLA from satellite observations (AVISO) and observationally-based reanalyses (GLORYS12, ECCO, and ORAS5) compared to the Month 1 hindcast from SEAS5, SPEAR, and CCSM4. The anomalies are relative to the climatology of 1995-2016 and relative to the global mean in January 1995. (a) Comparison of observed values with SEAS5 and SPEAR hindcasts, where the latter is corrected with its global mean steric component determined from the temperature and salinity profiles. (b) Comparison of Month 1 values from the SEAS5 and original (without global mean steric correction) SPEAR and CCSM4 hindcasts. (c) Comparison of linearly detrended ORAS5 with linearly trended

839 Month 1 values from the SEAS5, SPEAR, and CCSM4 hindcasts; the number in the legend  
 840 indicates the correlation of each hindcast time series with ORAS5.



841  
 842 Figure 4. Observed (gray) and predicted (red and blue) monthly SLA anomalies from August 1997 to  
 843 July 1998 at San Diego (top) and Charleston (bottom). Observations (dark gray) are from the tide  
 844 gauge records, and the model hindcasts were initialized by August 1, 1997 (LIM was initialized in  
 845 July 1997). The red solid line is the ensemble mean forecast, and the blue solid lines/  
 846 shading indicates ensemble members/spread. Units are cm.

847



848

849

850

851

852

853

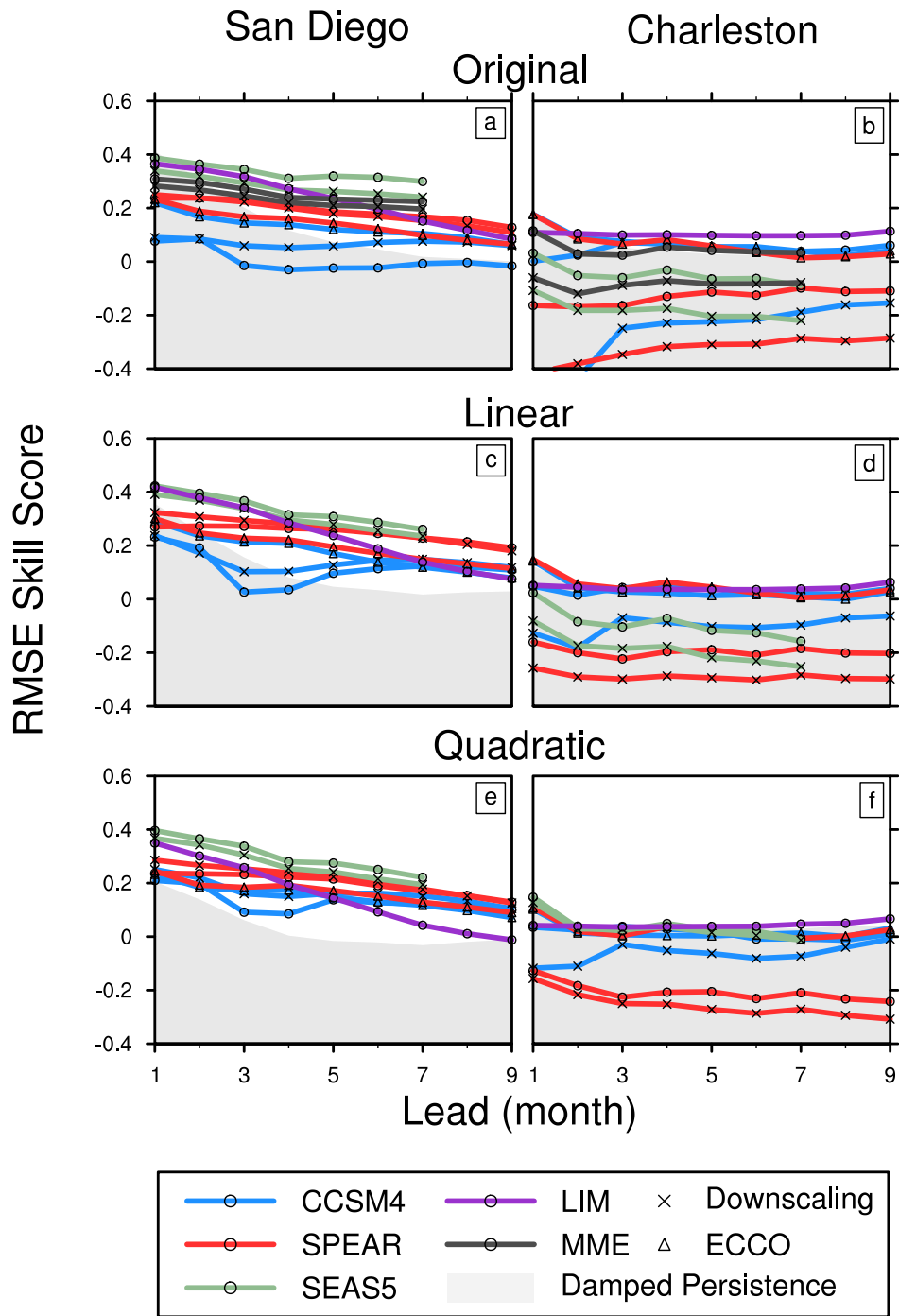
854

855

856

Figure 5. Deterministic skill measured by anomaly correlation coefficient (ACC) between the hindcasts and tide gauge observations at (left) San Diego and (right) Charleston at different lead times. The verification period is from 1995 to 2015, using hindcasts initialized in all calendar months. Gray shading shows damped persistence skill. Top row: Skill of each forecast technique for (a) San Diego and (b) Charleston. Second row: Same as (a and b) but after linear detrending of the observed tide gauge time series and from each hindcast. Third row: Same as (c and d) but using quadratic detrending. Note that trending also impacts the damped persistence time scale and skill.



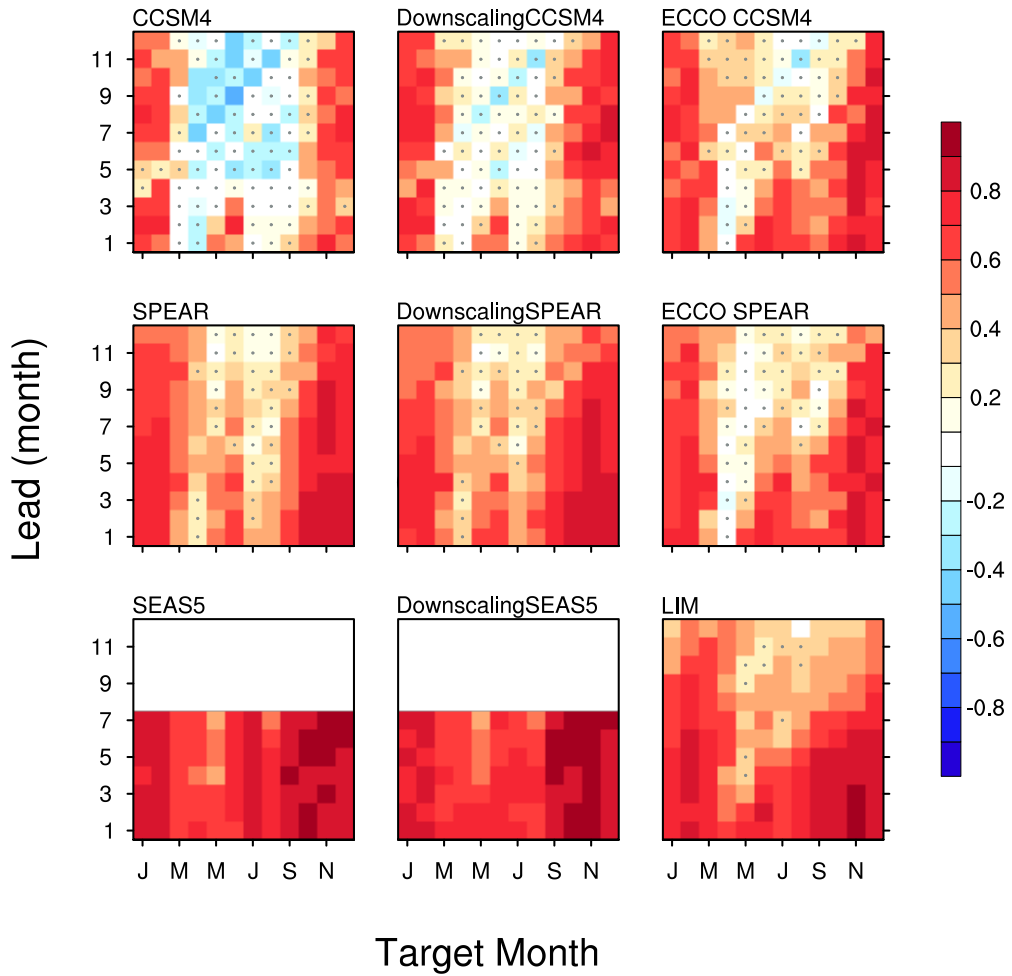


857

858 Figure 6. Same as Fig. 5 but using RMS skill score (RMSSS).

859

## San Diego

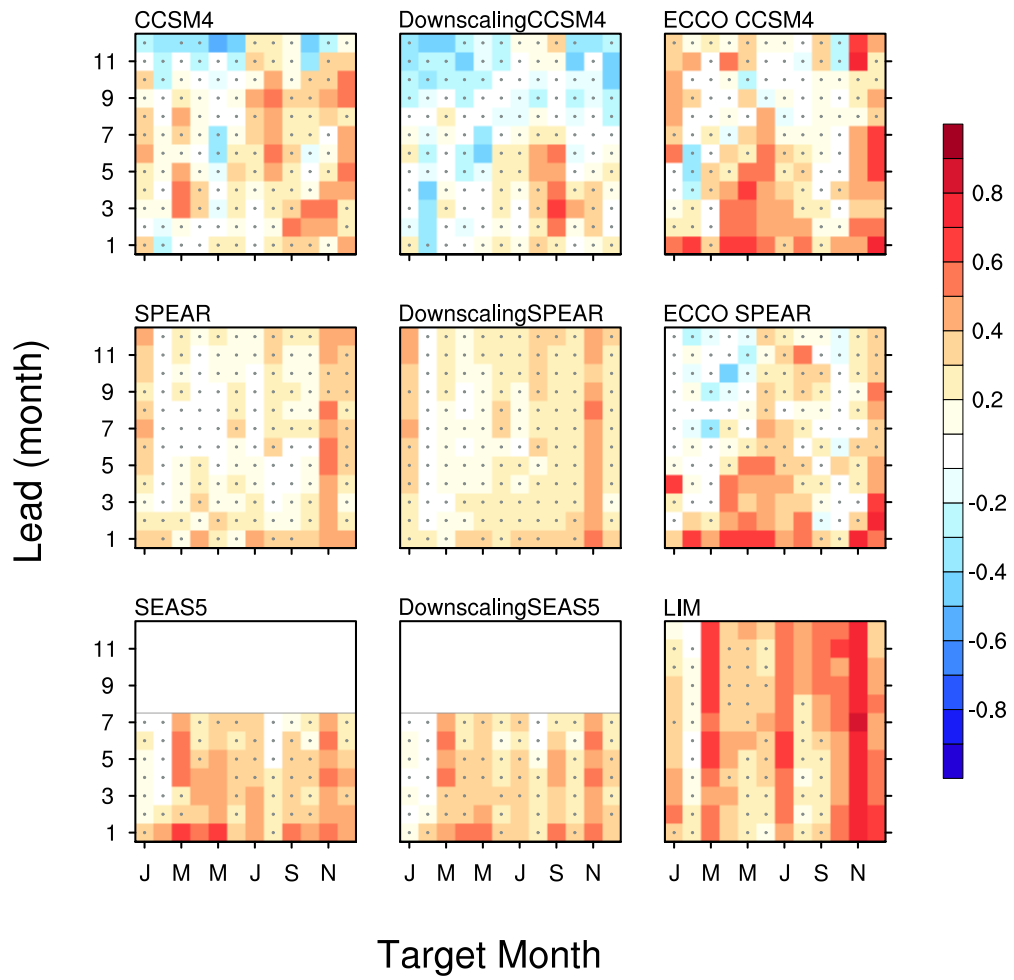


860

861 Figure 7. Deterministic skill measured by anomaly correlation coefficient (ACC) between the  
 862 hindcasts and tidal gauge observations at San Diego at different lead times and target months.  
 863 The gray dots indicate anomaly correlation values that are not significant at the 0.1 level  
 864 using a two-tail student t-test. The verification period is 1995 to 2015. No detrending is  
 865 performed upon either hindcasts or observations.

866

# Charleston

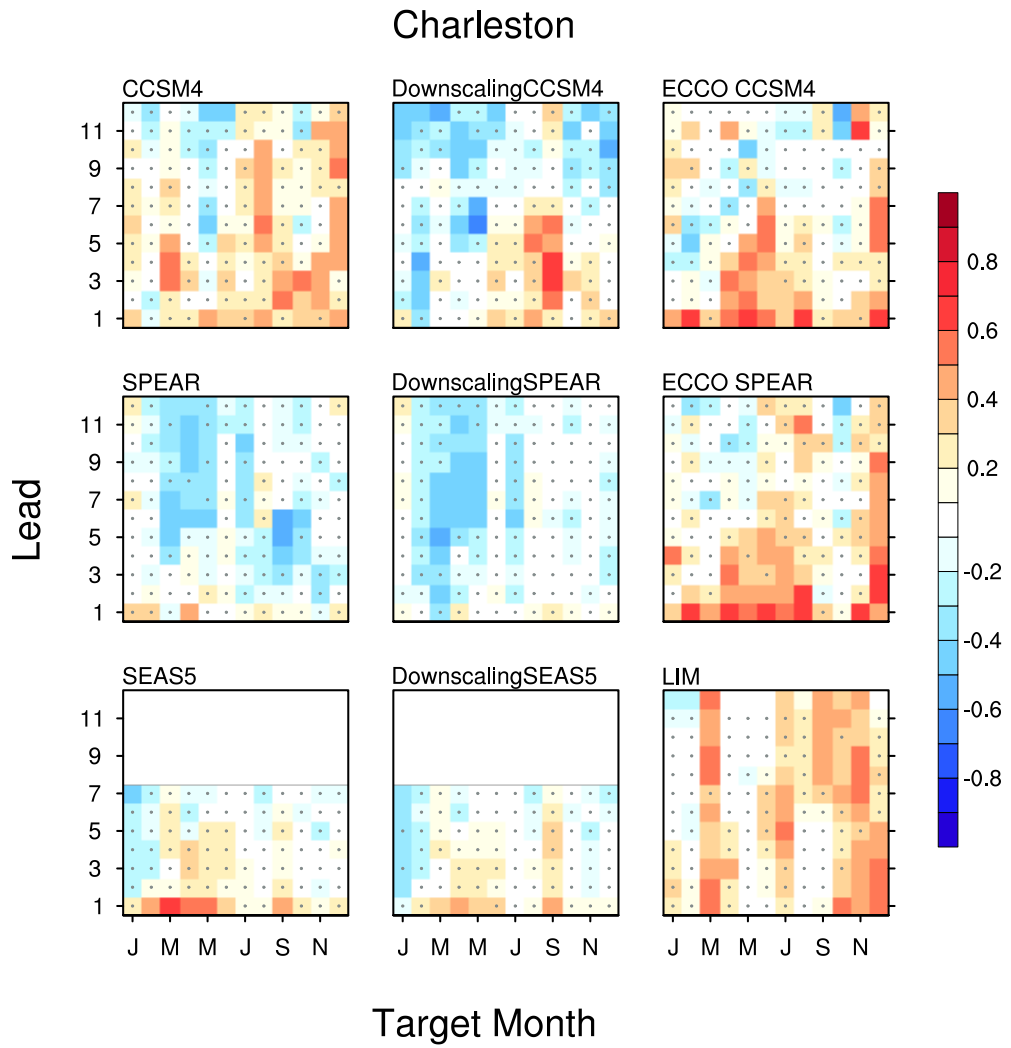


867

868

Figure 8. Same as Fig. 7 but for the Charleston tide gauge.

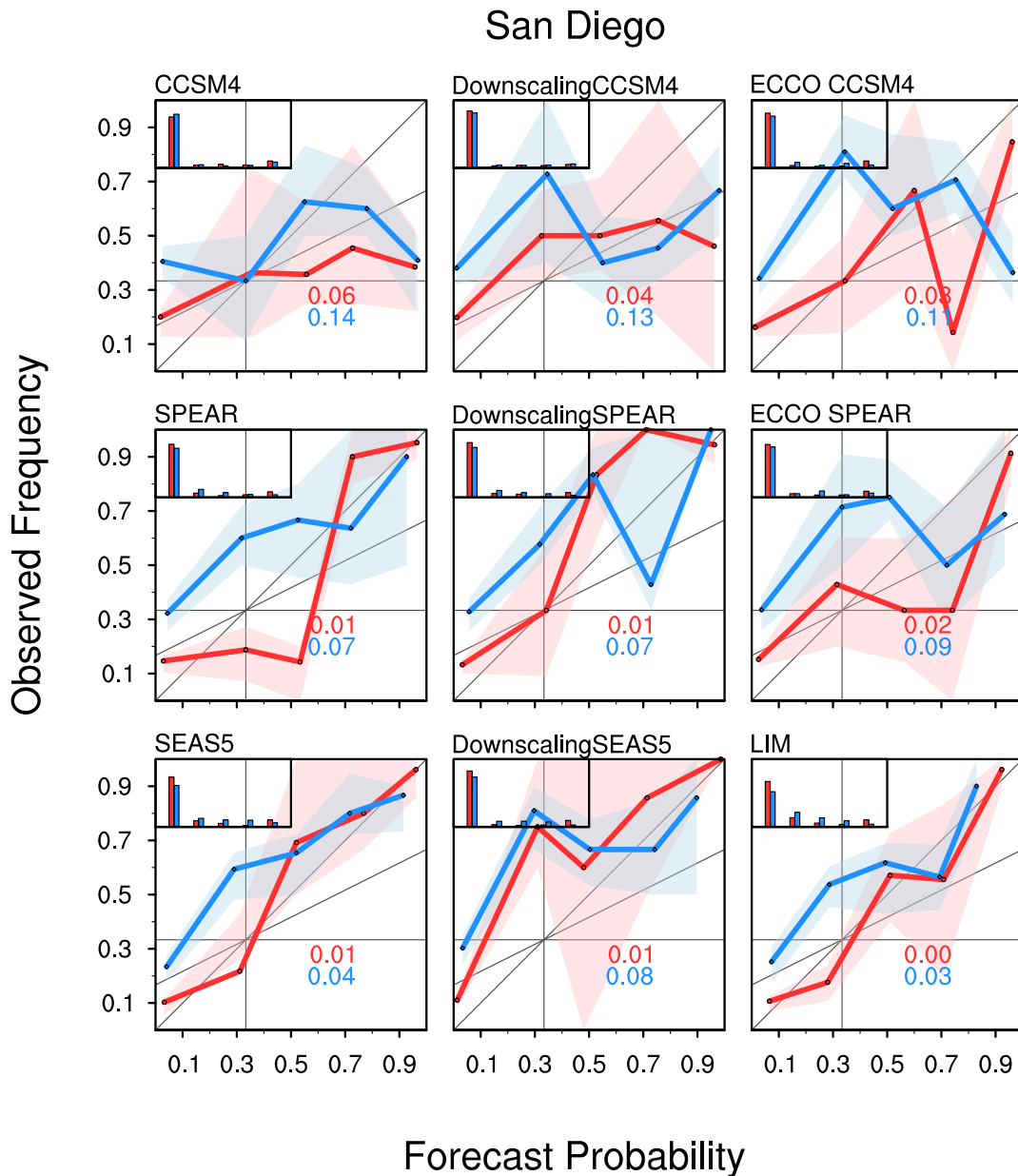
869



870

871 Figure 9. The same as Fig. 8 but computed after linearly detrending observations and  
 872 hindcasts.

873



874

875

876

877

878

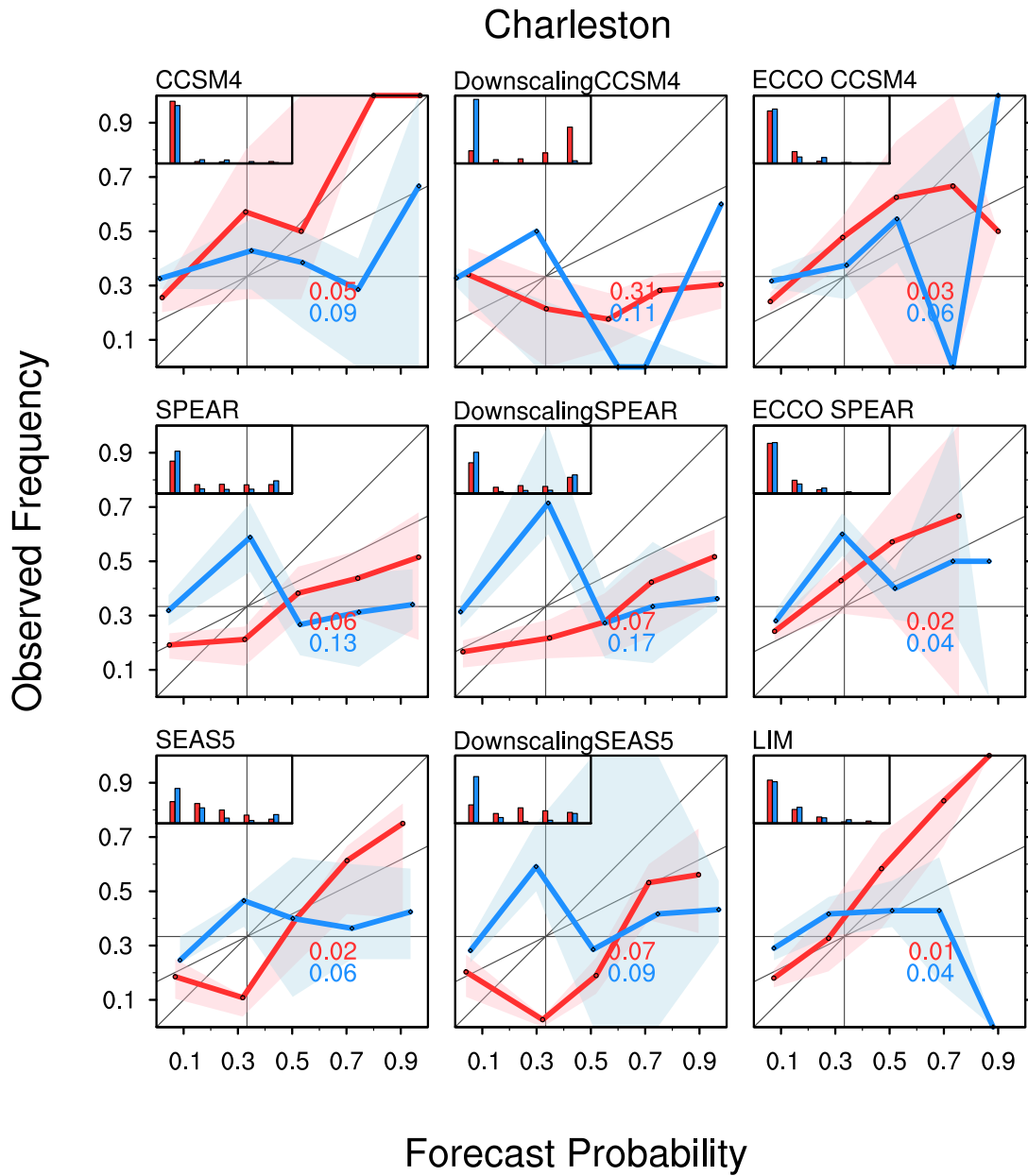
879

880

881

882

Figure 10. Reliability and sharpness diagrams of each hindcast for San Diego at Month 4. The mean forecast probability is plotted against the mean observed frequency for the reliability curve, determined by averaging all hindcasts within each quintile bin category. Red is for upper tercile hindcast, and blue is for lower tercile hindcast. The annotation is the reliability value with the same color coding (note that lower values represent better reliability). Gray shading shows the uncertainty of the reliability curves based on a bootstrapping calculation.

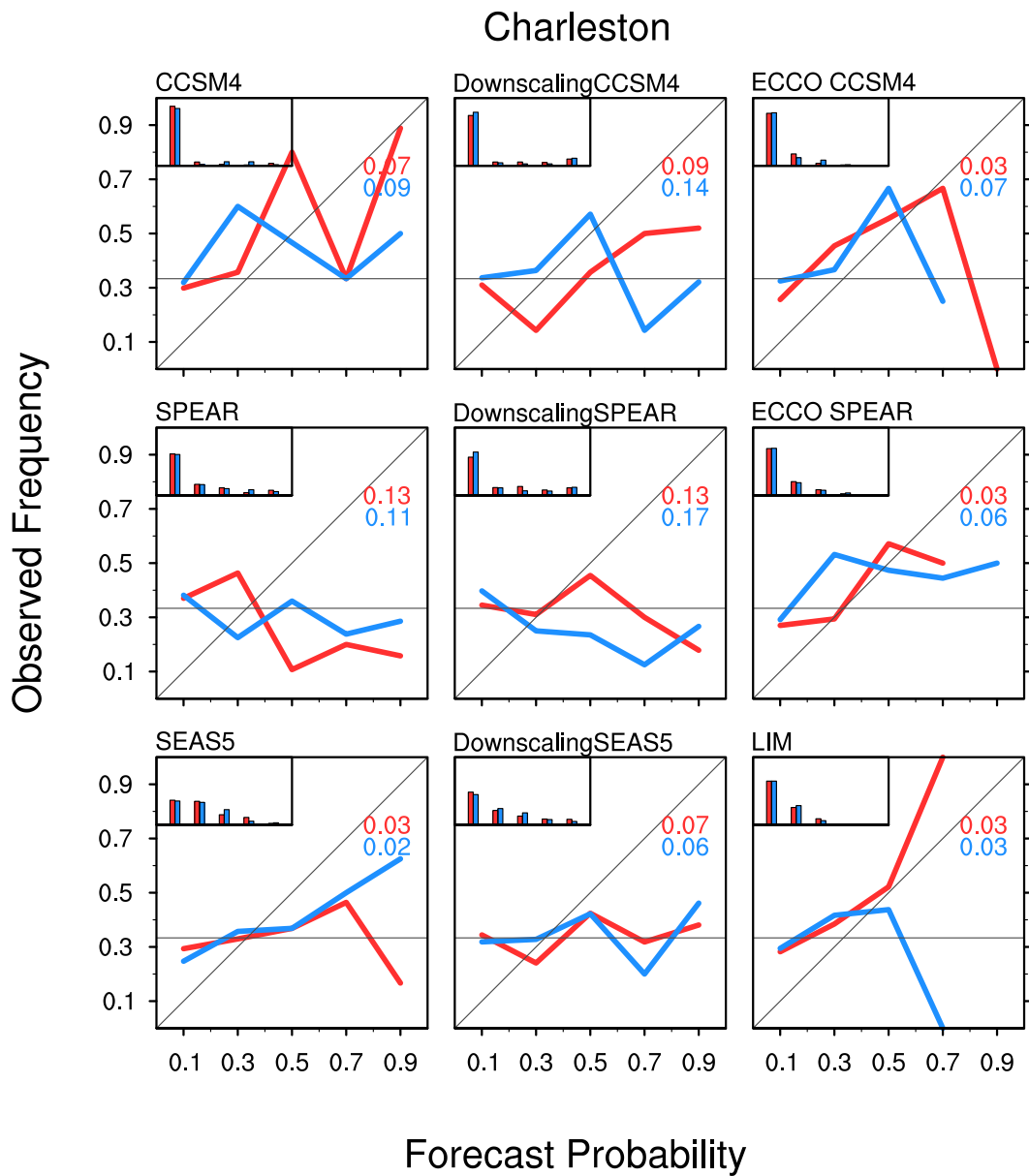


883

884

Figure 11. Same as Fig. 10 but for Charleston.

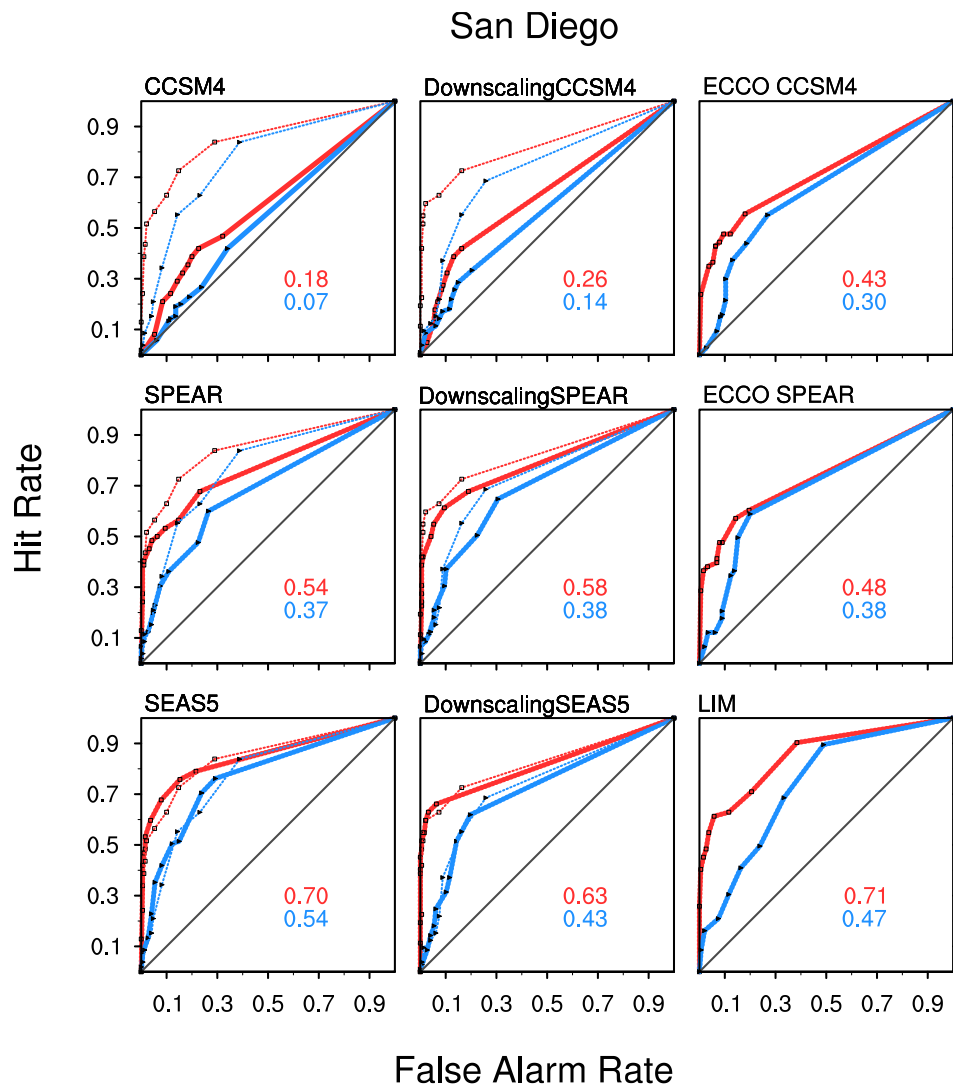
885



886

887 Figure 12. The same as Fig. 11 but for reliability and sharpness calculated after the linear  
 888 trend is removed from the observed tide gauge time series and each of the hindcasts. The  
 889 terciles of the observations are also calculated using the detrended data.

890

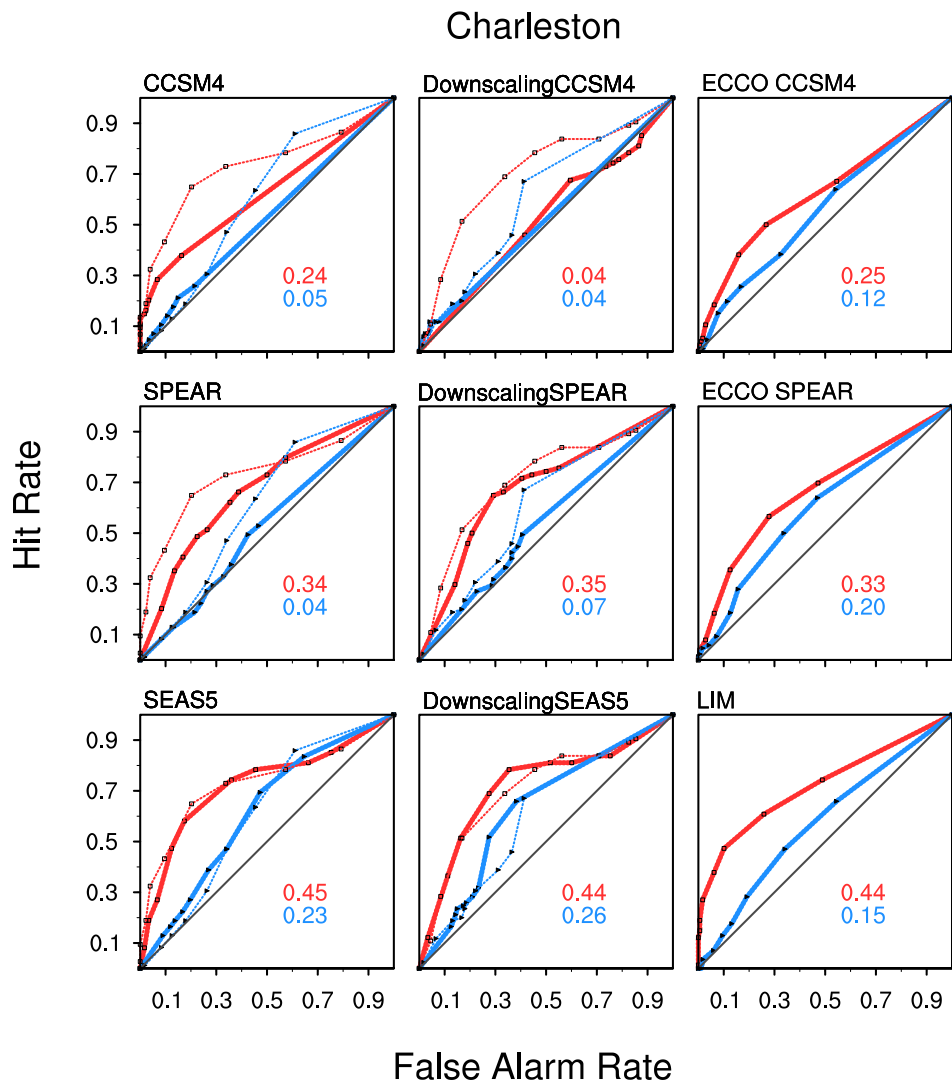


892

893 Figure 13. ROC curve of each of the hindcasts for San Diego at Month 4. Red is for the upper  
 894 tercile forecast, and blue is for the lower tercile forecast. The ROC skill score (ROCS) is  
 895 shown in each panel. The dashed lines in the first and second columns are the ROC curves for  
 896 the multi-model mean of the hindcasts from three dynamical models and three downscaled  
 897 versions, respectively.

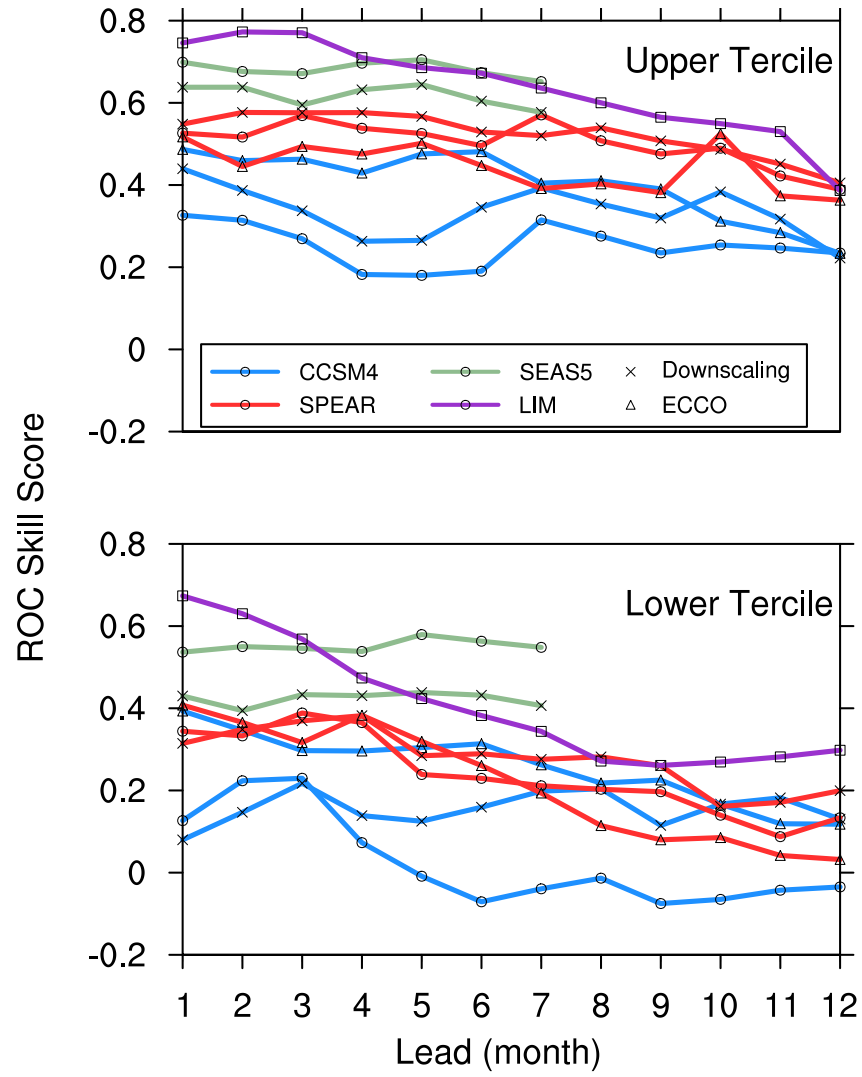
898





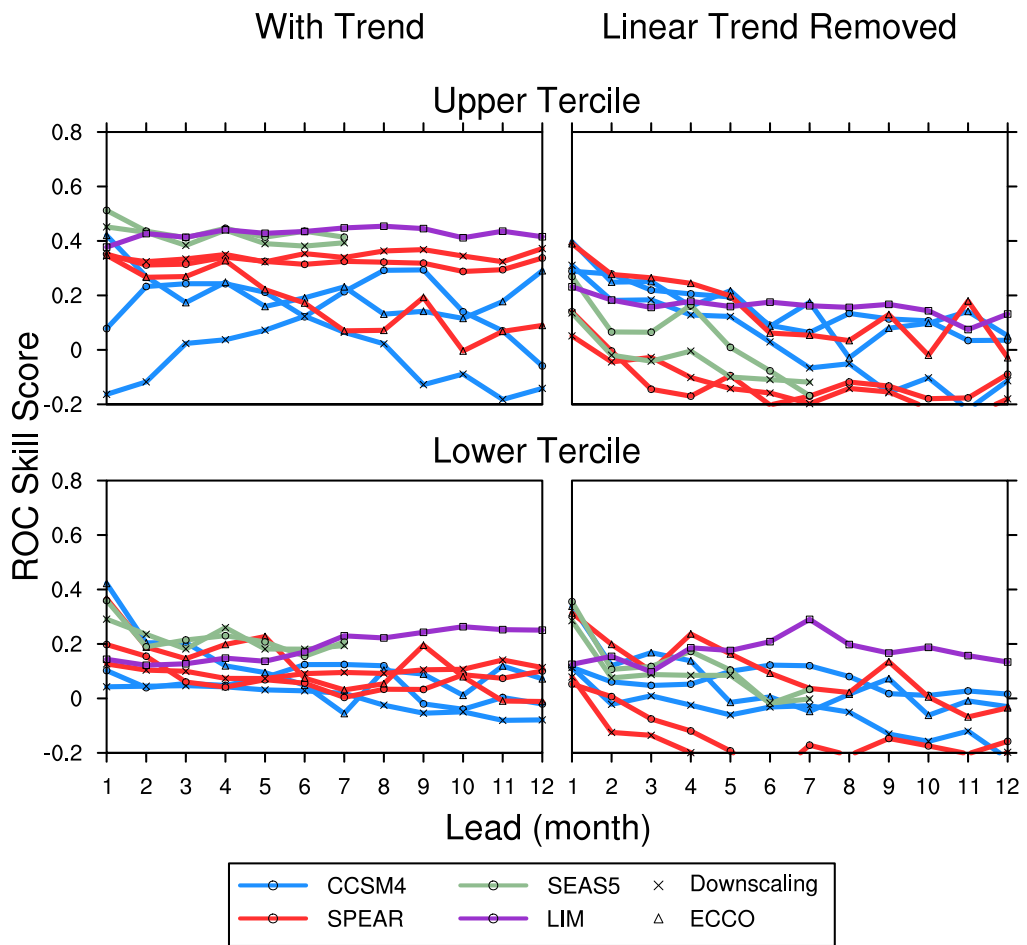
899  
 900  
 901

Figure 14. Same as Fig. 13 but for Charleston.



902

903 Figure 15. The ROC skill score (ROCS) for upper and lower tercile hindcast for each  
 904 prediction technique for San Diego at different lead times.



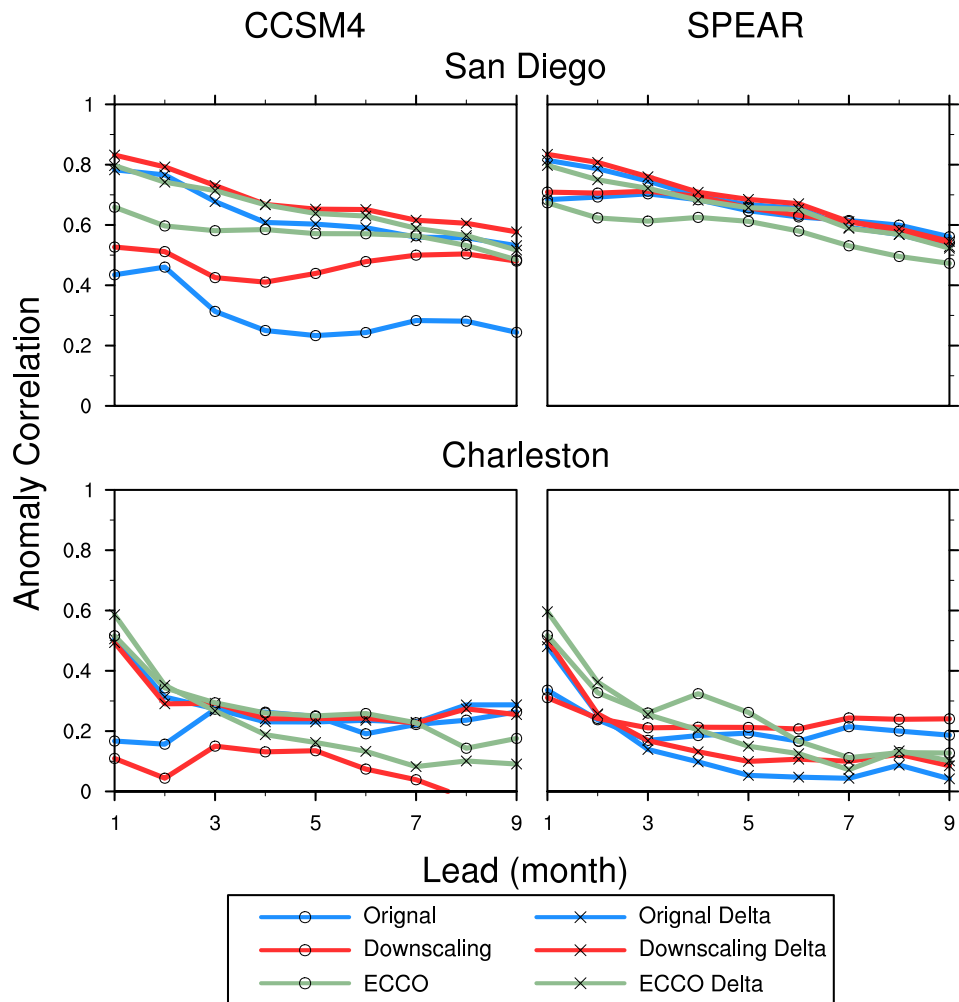
905

906 Figure 16. The ROC skill score (ROCS) for upper and lower tercile hindcast for each of the  
 907 prediction techniques, verified for Charleston at different lead times, determined for (left  
 908 column) full fields and (right column) linearly detrended fields.

909

910

911



912

913 Figure 17. Anomaly correlation coefficient (ACC) between the hindcasts and tidal gauge  
914 observations at San Diego and Charleston at different lead times. For each of the model  
915 hindcasts, the ACC of the delta corrected hindcast is compared with the skill of the original  
916 hindcast. See text for details of the procedure of creating the delta corrected hindcast.

917

918

919

920

- 922 Amaya DJ, MG Jacox, J Dias, MA Alexander, KB Karnauskas, JD Scott, and M Gehne  
 923 (2022). Subseasonal-to-seasonal forecast skill in the California Current System and its  
 924 connection to coastal Kelvin waves. *JGR: Oceans*.  
 925 <https://doi.org/10.1029/2021JC017892>.
- 926 Balmaseda, M. A., Mogensen, K., & Weaver, A. T. (2013). Evaluation of the ECMWF ocean  
 927 reanalysis system ORAS4. *Quarterly Journal of the Royal Meteorological Society*,  
 928 *139*(674), 1132-1161.
- 929 Balmaseda, M. A., Vitart, F., Ferranti, L., & Anderson, D. (2002). *Westerly wind events and*  
 930 *the 1997 El Nino event in the ECMWF seasonal forecasting system: A case study*.  
 931 ECMWF.
- 932 Beverley, J., Newman, M., & Hoell, A. (2023). Rapid Development of Systematic ENSO-  
 933 Related Seasonal Forecast Errors. *Geophysical Research Letters*, *50*(10),  
 934 e2022GL102249.
- 935 Brunner, K., D. Rivas, and K. M. M. Lwiza, 2019: Application of Classical Coastal Trapped  
 936 Wave Theory to High-Scattering Regions. *J. Phys. Oceanogr.*, **49**, 2201–2216
- 937 Bryan, F. O., Hecht, M. W., & Smith, R. D. (2007). Resolution convergence and sensitivity  
 938 studies with North Atlantic circulation models. Part I: The western boundary current  
 939 system. *Ocean Modelling*, *16*(3), 141-159.  
 940 <https://doi.org/https://doi.org/10.1016/j.ocemod.2006.08.005>
- 941 Cazenave, A., Henry, O., Munier, S., Delcroix, T., Gordon, A., Meyssignac, B., Llovel, W.,  
 942 Palanisamy, H., & Becker, M. (2012). Estimating ENSO influence on the global mean  
 943 sea level, 1993–2010. *Marine Geodesy*, *35*(sup1), 82-97.
- 944 Chowdhury, M. R., Chu, P. S., Schroeder, T., & Colasacco, N. (2007). Seasonal sea-level  
 945 forecasts by canonical correlation analysis—An operational scheme for the US-  
 946 affiliated Pacific Islands. *International Journal of Climatology: A Journal of the*  
 947 *Royal Meteorological Society*, *27*(10), 1389-1402.
- 948 Church, J. A., Clark, P. U., Cazenave, A., Gregory, J. M., Jevrejeva, S., Levermann, A.,  
 949 Merrifield, M. A., Milne, G. A., Nerem, R. S., & Nunn, P. D. (2013). Sea-level rise by  
 950 2100. *Science*, *342*(6165), 1445-1445.
- 951 Church, J. A., & White, N. J. (2006). A 20th century acceleration in global sea-level rise.  
 952 *Geophysical Research Letters*, *33*(1).
- 953 Consortium, E., Fukumori, I., Wang, O., Fenty, I., Forget, G., Heimbach, P., & Ponte, R.  
 954 (2020). ECCO central estimate (version 4 release 4). In.
- 955 Delworth, T. L., Cooke, W. F., Adcroft, A., Bushuk, M., Chen, J. H., Dunne, K. A., Ginoux,  
 956 P., Gudgel, R., Hallberg, R. W., & Harris, L. (2020). SPEAR: The next generation  
 957 GFDL modeling system for seasonal to multidecadal prediction and projection.  
 958 *Journal of Advances in Modeling Earth Systems*, *12*(3), e2019MS001895.
- 959 Deser, C., Phillips, A. S., Alexander, M. A., & Smoliak, B. V. (2014). Projecting North  
 960 American climate over the next 50 years: Uncertainty due to internal variability.  
 961 *Journal of Climate*, *27*(6), 2271-2296.
- 962 Ding, H., Newman, M., Alexander, M. A., & Wittenberg, A. T. (2018). Skillful climate  
 963 forecasts of the tropical Indo-Pacific Ocean using model-analogs. *Journal of Climate*,  
 964 *31*(14), 5437-5459.
- 965 Dusek, G., Sweet, W. V., Widlansky, M. J., Thompson, P. R., & Marra, J. J. (2022). A novel  
 966 statistical approach to predict seasonal high tide flooding [Original Research].  
 967 *Frontiers in Marine Science*, *9*. <https://doi.org/10.3389/fmars.2022.1073792>

- 968 Forget, G., Campin, J.-M., Heimbach, P., Hill, C., Ponte, R., & Wunsch, C. (2015). ECCO  
 969 version 4: An integrated framework for non-linear inverse modeling and global ocean  
 970 state estimation. *Geoscientific Model Development*, 8(10), 3071-3104.
- 971 Frankignoul, C., Gastineau, G., & Kwon, Y.-O. (2017). Estimation of the SST response to  
 972 anthropogenic and external forcing and its impact on the Atlantic multidecadal  
 973 oscillation and the Pacific decadal oscillation. *Journal of Climate*, 30(24), 9871-9895.
- 974 Frederikse, T., Lee, T., Wang, O., Kirtman, B., Becker, E., Hamlington, B., Limonadi, D., &  
 975 Waliser, D. (2022). A Hybrid Dynamical Approach for Seasonal Prediction of Sea-  
 976 Level Anomalies: A Pilot Study for Charleston, South Carolina. *Journal of*  
 977 *Geophysical Research: Oceans*, 127(8), e2021JC018137.  
 978 <https://doi.org/https://doi.org/10.1029/2021JC018137>
- 979 Garcia-Soto, C., Cheng, L., Caesar, L., Schmidtko, S., Jewett, E. B., Cheripka, A., Rigor, I.,  
 980 Caballero, A., Chiba, S., & Báez, J. C. (2021). An overview of ocean climate change  
 981 indicators: Sea surface temperature, ocean heat content, ocean pH, dissolved oxygen  
 982 concentration, arctic sea ice extent, thickness and volume, sea level and strength of  
 983 the AMOC (Atlantic Meridional Overturning Circulation). *Frontiers in Marine*  
 984 *Science*, 8, 642372.
- 985 Gill, S. K., & Schultz, J. R. (2001). Tidal datums and their applications.
- 986 Griffies, S. M., & Adcroft, A. J. (2008). Formulating the equations of ocean models. *Ocean*  
 987 *modeling in an eddying regime*, 177, 281-317.
- 988 Griffies, S. M., Danabasoglu, G., Durack, P. J., Adcroft, A. J., Balaji, V., Böning, C. W.,  
 989 Chassignet, E. P., Curchitser, E., Deshayes, J., & Drange, H. (2016). OMIP  
 990 contribution to CMIP6: Experimental and diagnostic protocol for the physical  
 991 component of the Ocean Model Intercomparison Project. *Geoscientific Model*  
 992 *Development*, 3231.
- 993 Griffies, S. M., & Greatbatch, R. J. (2012). Physical processes that impact the evolution of  
 994 global mean sea level in ocean climate models. *Ocean Modelling*, 51, 37-72.  
 995 <https://doi.org/https://doi.org/10.1016/j.ocemod.2012.04.003>
- 996 Hague, B. S., McGregor, S., Jones, D. A., Reef, R., Jakob, D., & Murphy, B. F. (2023). The  
 997 Global Drivers of Chronic Coastal Flood Hazards Under Sea-Level Rise. *Earth's*  
 998 *Future*, 11(8), e2023EF003784. <https://doi.org/https://doi.org/10.1029/2023EF003784>
- 999 Hamlington, B. D., Gardner, A. S., Ivins, E., Lenaerts, J. T., Reager, J., Trossman, D. S.,  
 1000 Zaron, E. D., Adhikari, S., Arendt, A., & Aschwanden, A. (2020). Understanding of  
 1001 contemporary regional sea-level change and the implications for the future. *Reviews*  
 1002 *of Geophysics*, 58(3), e2019RG000672.
- 1003 Hamlington, B. D., Leben, R. R., Kim, K.-Y., Nerem, R. S., Atkinson, L. P., & Thompson, P.  
 1004 R. (2015). The effect of the El Niño-Southern Oscillation on U.S. regional and coastal  
 1005 sea level. *Journal of Geophysical Research: Oceans*, 120(6), 3970-3986.  
 1006 <https://doi.org/https://doi.org/10.1002/2014JC010602>
- 1007 Hamlington, B. D., Picuch, C. G., Reager, J. T., Chandanpurkar, H., Frederikse, T., Nerem,  
 1008 R. S., Fasullo, J. T., & Cheon, S.-H. (2020). Origin of interannual variability in global  
 1009 mean sea level. *Proceedings of the National Academy of Sciences*, 117(25), 13983-  
 1010 13990. <https://doi.org/doi:10.1073/pnas.1922190117>
- 1011 Han, W., Meehl, G. A., Stammer, D., Hu, A., Hamlington, B., Kenigson, J., Palanisamy, H.,  
 1012 & Thompson, P. (2017). Spatial patterns of sea level variability associated with  
 1013 natural internal climate modes. *Integrative study of the mean sea level and its*  
 1014 *components*, 221-254.
- 1015 Han, W., Stammer, D., Thompson, P., Ezer, T., Palanisamy, H., Zhang, X., Domingues, C.  
 1016 M., Zhang, L., & Yuan, D. (2019a). Impacts of Basin-Scale Climate Modes on

1017 Coastal Sea Level: a Review. *Surveys in geophysics*, 40(6), 1493-1541.  
1018 <https://doi.org/10.1007/s10712-019-09562-8>

1019 Han, W., Stammer, D., Thompson, P., Ezer, T., Palanisamy, H., Zhang, X., Domingues, C.  
1020 M., Zhang, L., & Yuan, D. (2019b). Impacts of basin-scale climate modes on coastal  
1021 sea level: a review. *Surveys in geophysics*, 40, 1493-1541.

1022 Hauser, D., Tourain, C., Hermozo, L., Alraddawi, D., Aouf, L., Chapron, B., Dalphiné, A.,  
1023 Delaye, L., Dalila, M., Dormy, E., Gouillon, F., Gressani, V., Grouazel, A., Guitton,  
1024 G., Husson, R., Mironov, A., Mouche, A., Ollivier, A., Oruba, L., . . . Tran, N. (2021).  
1025 New Observations From the SWIM Radar On-Board CFOSAT: Instrument Validation  
1026 and Ocean Wave Measurement Assessment. *IEEE Transactions on Geoscience and*  
1027 *Remote Sensing*, 59(1), 5-26. <https://doi.org/10.1109/TGRS.2020.2994372>

1028 Holgate, S. J., Matthews, A., Woodworth, P. L., Rickards, L. J., Tamisiea, M. E., Bradshaw,  
1029 E., Foden, P. R., Gordon, K. M., Jevrejeva, S., & Pugh, J. (2013). New data systems  
1030 and products at the permanent service for mean sea level. *Journal of Coastal*  
1031 *Research*, 29(3), 493-504.

1032 Hughes, C.W., Fukumori, I., Griffies, S.M. *et al.*, 2019: Sea Level and the Role of Coastal  
1033 Trapped Waves in Mediating the Influence of the Open Ocean on the Coast. *Surv*  
1034 *Geophys* 40, 1467–1492 .

1035 Jean-Michel, L., Eric, G., Romain, B.-B., Gilles, G., Angélique, M., Marie, D., Clément, B.,  
1036 Mathieu, H., Olivier, L. G., & Charly, R. (2021). The Copernicus global 1/12 oceanic  
1037 and sea ice GLORYS12 reanalysis. *Frontiers in Earth Science*, 9, 698876.

1038 Johnson, S. J., Stockdale, T. N., Ferranti, L., Balmaseda, M. A., Molteni, F., Magnusson, L.,  
1039 Tietsche, S., Decremer, D., Weisheimer, A., & Balsamo, G. (2019). SEAS5: the new  
1040 ECMWF seasonal forecast system. *Geoscientific Model Development*, 12(3), 1087-  
1041 1117.

1042 Johnson, S. J., Stockdale, T. N., Ferranti, L., Balmaseda, M. A., Molteni, F., Magnusson, L.,  
1043 Tietsche, S., Decremer, D., Weisheimer, A., Balsamo, G., Keeley, S. P. E., Mogensen,  
1044 K., Zuo, H., & Monge-Sanz, B. M. (2019). SEAS5: the new ECMWF seasonal  
1045 forecast system. *Geosci. Model Dev.*, 12(3), 1087-1117. <https://doi.org/10.5194/gmd-12-1087-2019>

1046 Kennedy, J. J., Rayner, N. A., Atkinson, C. P., & Killick, R. E. (2019). An Ensemble Data  
1047 Set of Sea Surface Temperature Change From 1850: The Met Office Hadley Centre  
1048 HadSST.4.0.0.0 Data Set. *Journal of Geophysical Research: Atmospheres*, 124(14),  
1049 7719-7763. <https://doi.org/https://doi.org/10.1029/2018JD029867>

1050 Kharin, V. V., & Zwiers, F. W. (2003). On the ROC score of probability forecasts. *Journal of*  
1051 *Climate*, 16(24), 4145-4150.

1052 Kirtman, B. P., Min, D., Infanti, J. M., Kinter, J. L., Paolino, D. A., Zhang, Q., Van Den  
1053 Dool, H., Saha, S., Mendez, M. P., & Becker, E. (2014). The North American  
1054 multimodel ensemble: phase-1 seasonal-to-interannual prediction; phase-2 toward  
1055 developing intraseasonal prediction. *Bulletin of the American Meteorological Society*,  
1056 95(4), 585-601.

1057 Kirtman, B. P., Min, D., Infanti, J. M., Kinter, J. L., Paolino, D. A., Zhang, Q., van den Dool,  
1058 H., Saha, S., Mendez, M. P., Becker, E., Peng, P., Tripp, P., Huang, J., DeWitt, D. G.,  
1059 Tippett, M. K., Barnston, A. G., Li, S., Rosati, A., Schubert, S. D., . . . Wood, E. F.  
1060 (2014). The North American Multimodel Ensemble: Phase-1 Seasonal-to-Interannual  
1061 Prediction; Phase-2 toward Developing Intraseasonal Prediction. *Bulletin of the*  
1062 *American Meteorological Society*, 95(4), 585-601.  
1063 <https://doi.org/https://doi.org/10.1175/BAMS-D-12-00050.1>

1064 Livezey, R. E., & Chen, W. (1983). Statistical field significance and its determination by  
1065 Monte Carlo techniques. *Mon. Wea. Rev.*, 111(1), 46-59.



- 1067 Long, X., Shin, S. I., & Newman, M. (2023). Statistical downscaling of seasonal forecasts of  
 1068 sea level anomalies for US coasts. *Geophysical Research Letters*, 50(4),  
 1069 e2022GL100271.
- 1070 Long, X., Widlansky, M. J., Schloesser, F., Thompson, P. R., Annamalai, H., Merrifield, M.  
 1071 A., & Yoon, H. (2020). Higher Sea Levels at Hawaii Caused by Strong El Niño and  
 1072 Weak Trade Winds. *Journal of Climate*, 33(8), 3037-3059.  
 1073 <https://doi.org/https://doi.org/10.1175/JCLI-D-19-0221.1>
- 1074 Long, X., Widlansky, M. J., Spillman, C. M., Kumar, A., Balmaseda, M., Thompson, P. R.,  
 1075 Chikamoto, Y., Smith, G. A., Huang, B., Shin, C.-S., Merrifield, M. A., Sweet, W. V.,  
 1076 Leuliette, E., Annamalai, H. S., Marra, J. J., & Mitchum, G. (2021). Seasonal  
 1077 Forecasting Skill of Sea-Level Anomalies in a Multi-Model Prediction Framework.  
 1078 *Journal of Geophysical Research: Oceans*, 126(6), e2020JC017060.  
 1079 <https://doi.org/https://doi.org/10.1029/2020JC017060>
- 1080 Lu, F., Harrison, M. J., Rosati, A., Delworth, T. L., Yang, X., Cooke, W. F., Jia, L., McHugh,  
 1081 C., Johnson, N. C., Bushuk, M., Zhang, Y., & Adcroft, A. (2020). GFDL's SPEAR  
 1082 Seasonal Prediction System: Initialization and Ocean Tendency Adjustment (OTA)  
 1083 for Coupled Model Predictions. *Journal of Advances in Modeling Earth Systems*,  
 1084 12(12), e2020MS002149. <https://doi.org/https://doi.org/10.1029/2020MS002149>
- 1085 Mason, S. J., & Graham, N. E. (1999). Conditional probabilities, relative operating  
 1086 characteristics, and relative operating levels. *Weather and Forecasting*, 14(5), 713-  
 1087 725.
- 1088 May, C. L., Osler, M. S., Stockdon, H. F., Barnard, P. L., Callahan, J. A., Collini, R. C.,  
 1089 Ferreira, C. M., Finzi Hart, J., Lentz, E. E., Mahoney, T. B., Sweet, W., Walker, D., &  
 1090 Weaver, C. P. (2023). Coastal effects. In A. R. Crimmins, C. W. Avery, D. R.  
 1091 Easterling, K. E. Kunkel, B. C. Stewart, & T. K. Maycock (Eds.), *Fifth National*  
 1092 *Climate Assessment*. U.S. Global Change Research Program.  
 1093 <https://doi.org/10.7930/NCA5.2023.CH9>
- 1094 McIntosh, P. C., Church, J. A., Miles, E. R., Ridgway, K., & Spillman, C. M. (2015).  
 1095 Seasonal coastal sea level prediction using a dynamical model. *Geophysical Research*  
 1096 *Letters*, 42(16), 6747-6753.
- 1097 Miles, E. R., Spillman, C. M., Church, J. A., & McIntosh, P. C. (2014). Seasonal prediction  
 1098 of global sea level anomalies using an ocean–atmosphere dynamical model. *Climate*  
 1099 *dynamics*, 43, 2131-2145.
- 1100 Newman, M., & Sardeshmukh, P. D. (2017). Are we near the predictability limit of tropical  
 1101 Indo-Pacific sea surface temperatures? *Geophysical Research Letters*, 44(16), 8520-  
 1102 8529.
- 1103 Nicholls, R. J., Lincke, D., Hinkel, J., Brown, S., Vafeidis, A. T., Meyssignac, B., Hanson, S.  
 1104 E., Merkens, J.-L., & Fang, J. (2021). A global analysis of subsidence, relative sea-  
 1105 level change and coastal flood exposure. *Nature Climate Change*, 11(4), 338-342.  
 1106 <https://doi.org/10.1038/s41558-021-00993-z>
- 1107 Parker, B. B. (2007). Tidal analysis and prediction.
- 1108 Penland, C., & Sardeshmukh, P. D. (1995). The Optimal Growth of Tropical Sea Surface  
 1109 Temperature Anomalies. *Journal of Climate*, 8(8), 1999-2024.  
 1110 [https://doi.org/https://doi.org/10.1175/1520-  
 1111 0442\(1995\)008<1999:TOGOTS>2.0.CO;2](https://doi.org/https://doi.org/10.1175/1520-0442(1995)008<1999:TOGOTS>2.0.CO;2)
- 1112 Piecuch, C. G., Bittermann, K., Kemp, A. C., Ponte, R. M., Little, C. M., Engelhart, S. E., &  
 1113 Lentz, S. J. (2018). River-discharge effects on United States Atlantic and Gulf coast  
 1114 sea-level changes. *Proceedings of the National Academy of Sciences*, 115(30), 7729-  
 1115 7734.



- 1116 Pugh, D., & Woodworth, P. (2014). *Sea-level science: understanding tides, surges, tsunamis*  
1117 *and mean sea-level changes*. Cambridge University Press.
- 1118 Ray, R., Widlansky, M., Genz, A., & Thompson, P. (2023). Offsets in tide-gauge reference  
1119 levels detected by satellite altimetry: ten case studies. *Journal of Geodesy*, 97(12),  
1120 110.
- 1121 Roberts, C., Calvert, D., Dunstone, N., Hermanson, L., Palmer, M., & Smith, D. (2016). On  
1122 the drivers and predictability of seasonal-to-interannual variations in regional sea  
1123 level. *Journal of Climate*, 29(21), 7565-7585.
- 1124 Ryan, H., & Noble, M. (2002). Sea level response to ENSO along the central California  
1125 coast: how the 1997–1998 event compares with the historic record. *Progress in*  
1126 *Oceanography*, 54(1-4), 149-169.
- 1127 Sardeshmukh, P. D., J. A. Wang, G. P. Compo, and C. Penland, 2023: Improving  
1128 Atmospheric Models by Accounting for Chaotic Physics. *J. Climate*, 36, 5569–5585
- 1129 Shin, S.-I., & Newman, M. (2021). Seasonal Predictability of Global and North American  
1130 Coastal Sea Surface Temperature and Height Anomalies. *Geophysical Research*  
1131 *Letters*, 48(10), e2020GL091886.  
1132 <https://doi.org/https://doi.org/10.1029/2020GL091886>
- 1133 Smith, D. M., Eade, R., & Pohlmann, H. (2013). A comparison of full-field and anomaly  
1134 initialization for seasonal to decadal climate prediction. *Climate dynamics*, 41(11),  
1135 3325-3338. <https://doi.org/10.1007/s00382-013-1683-2>
- 1136 Taherkhani, M., Vitousek, S., Barnard, P. L., Frazer, N., Anderson, T. R., & Fletcher, C. H.  
1137 (2020). Sea-level rise exponentially increases coastal flood frequency. *Scientific*  
1138 *reports*, 10(1), 1-17.
- 1139 Thompson, P. R., Widlansky, M. J., Hamlington, B. D., Merrifield, M. A., Marra, J. J.,  
1140 Mitchum, G. T., & Sweet, W. (2021). Rapid increases and extreme months in  
1141 projections of United States high-tide flooding. *Nature Climate Change*, 11(7), 584-  
1142 590. <https://doi.org/10.1038/s41558-021-01077-8>
- 1143 Tippett, M. K., & L'Heureux, M. L. (2020). Low-dimensional representations of Niño 3.4  
1144 evolution and the spring persistence barrier. *npj Climate and Atmospheric Science*,  
1145 3(1), 24.
- 1146 Toth, Z., Talagrand, O., & Zhu, Y. (2006). The attributes of forecast systems: A general  
1147 framework for the evaluation and calibration of weather forecasts. *Predictability of*  
1148 *Weather and Climate*, 584(2006), 595.
- 1149 van den Dool, H. (2006). *Empirical Methods in Short-Term Climate Prediction*. Oxford  
1150 University Press. <https://doi.org/10.1093/oso/9780199202782.001.0001>
- 1151 Vitousek, S., Barnard, P. L., Fletcher, C. H., Frazer, N., Erikson, L., & Storlazzi, C. D.  
1152 (2017). Doubling of coastal flooding frequency within decades due to sea-level rise.  
1153 *Scientific reports*, 7(1), 1-9.
- 1154 Wang, G., Ren, H.-L., Liu, J., & Long, X. (2023). Seasonal predictions of sea surface height  
1155 in BCC-CSM1.1m and their modulation by tropical climate dominant modes.  
1156 *Atmospheric Research*, 281, 106466.  
1157 <https://doi.org/https://doi.org/10.1016/j.atmosres.2022.106466>
- 1158 Wang, J., Church, J. A., Zhang, X., & Chen, X. (2021). Reconciling global mean and regional  
1159 sea level change in projections and observations. *Nature Communications*, 12(1), 990.  
1160 <https://doi.org/10.1038/s41467-021-21265-6>
- 1161 Weisheimer, A., & Palmer, T. N. (2014). On the reliability of seasonal climate forecasts.  
1162 *Journal of the Royal Society Interface*, 11(96), 20131162.
- 1163 Widlansky, M. J., Long, X., Balmaseda, M. A., Spillman, C. M., Smith, G., Zuo, H., Yin, Y.,  
1164 Alves, O., & Kumar, A. (2023). Quantifying the Benefits of Altimetry Assimilation in

1165 Seasonal Forecasts of the Upper Ocean. *Journal of Geophysical Research: Oceans*,  
1166 128(5), e2022JC019342. <https://doi.org/https://doi.org/10.1029/2022JC019342>

1167 Widlansky, M. J., Marra, J. J., Chowdhury, M. R., Stephens, S. A., Miles, E. R., Fauchereau,  
1168 N., Spillman, C. M., Smith, G., Beard, G., & Wells, J. (2017). Multimodel ensemble  
1169 sea level forecasts for tropical Pacific islands. *Journal of Applied Meteorology and*  
1170 *Climatology*, 56(4), 849-862.

1171 Wilks, D. S. (2011). *Statistical methods in the atmospheric sciences* (Vol. 100). Academic  
1172 press.

1173 Wöppelmann, G., & Marcos, M. (2016). Vertical land motion as a key to understanding sea  
1174 level change and variability. *Reviews of Geophysics*, 54(1), 64-92.  
1175 <https://doi.org/https://doi.org/10.1002/2015RG000502>

1176 Wulff, C. O., Vitart, F., & Domeisen, D. I. V. (2022). Influence of trends on subseasonal  
1177 temperature prediction skill. *Quarterly Journal of the Royal Meteorological Society*,  
1178 148(744), 1280-1299. <https://doi.org/https://doi.org/10.1002/qj.4259>

1179 Xie, D., Zou, Q.-P., Mignone, A., & MacRae, J. D. (2019). Coastal flooding from wave  
1180 overtopping and sea level rise adaptation in the northeastern USA. *Coastal*  
1181 *Engineering*, 150, 39-58.

1182 Xue, Y., Huang, B., Hu, Z.-Z., Kumar, A., Wen, C., Behringer, D., & Nadiga, S. (2011). An  
1183 assessment of oceanic variability in the NCEP climate forecast system reanalysis.  
1184 *Climate dynamics*, 37(11), 2511-2539. <https://doi.org/10.1007/s00382-010-0954-4>

1185 Zervas, C. E., Gill, S. K., & Sweet, W. W. V. (2013). Estimating vertical land motion from  
1186 long-term tide gauge records.

1187 Zuo, H., Balmaseda, M. A., Tietsche, S., Mogensen, K., & Mayer, M. (2019). The ECMWF  
1188 operational ensemble reanalysis–analysis system for ocean and sea ice: a description  
1189 of the system and assessment. *Ocean science*, 15(3), 779-808.

1190