

# Assessing the storm surge model performance: What error indicators can measure the skill?

Rodrigo Campos-Caba<sup>1</sup>, Lorenzo Mentaschi<sup>1</sup>, Jacopo Alessandri<sup>1</sup>, Paula Camus<sup>2</sup>, Andrea Mazzino<sup>3</sup>, Francesco Ferrari<sup>3</sup>, Ivan Federico<sup>4</sup>, Michalis Vousdoukas<sup>5</sup>, and Massimo Tondello<sup>6</sup>

5 <sup>1</sup>Department of Physics and Astronomy, University of Bologna, Bologna, Italy

<sup>2</sup>Departamento de Ciencias y Técnicas del Agua y del Medio Ambiente, University of Cantabria, Santander, Spain

<sup>3</sup>Department of Civil, Chemical and Environmental Engineering, University of Genoa, Genoa, Italy

<sup>4</sup>Centro Euro-Mediterraneo sui Cambiamenti Climatici – Ocean Predictions and Applications, Lecce, Italy

<sup>5</sup>Joint Research Centre, European Commission, Seville, Spain

10 <sup>6</sup>HS Marine Srl, Noventa Padovana, Italy

*Correspondence to:* Rodrigo Campos-Caba (rodrigo.camposcaba@unibo.it)

**Abstract.** A well-validated storm surge numerical model is crucial, offering precise coastal hazard information and serving as a basis for extensive databases and advanced data-driven algorithms. However, selecting the best model setup based solely on common error indicators like RMSE or Pearson correlation doesn't always yield optimal results. To illustrate this, we conducted 34-year high-resolution simulations for storm surge under barotropic (BT) and baroclinic (BC) configurations, using atmospheric data from ERA5 and a high-resolution downscaling of the Climate Forecast System Reanalysis (CFSR) developed by the University of Genoa (UniGe). We combined forcing and configurations to produce three datasets: 1) BT-ERA5, 2) BC-ERA5, and 3) BC-UniGe. The model performance was assessed against nearshore station data using various statistical metrics. While RMSE and Pearson correlation suggest BT-ERA5, i.e. the coarsest and simplest setup, as the best model, followed by BC-ERA5, we demonstrate that these indicators aren't always reliable for performance assessment. The most sophisticated model BC-UniGe shows worse values of RMSE or Pearson correlation due to the so-called “double penalty” effect. Here we propose new skill indicators that assess the ability of the model to reproduce the distribution of the observations. This, combined with an analysis of values above the 99th percentile, identifies BC-UniGe as the best model, while ERA5 simulations tend to underestimate the extremes. Although the study focuses on the accurate representation of storm surge by the numerical model, the analysis and proposed metrics can be applied to any problem involving the comparison between time series of simulation and observation.

15  
20  
25

30

## 1 Introduction

35 In coastal areas, accurately depicting storm surge is paramount for effective risk assessment, preparedness, and mitigation strategies, as they can lead to coastal erosion, inundation, and infrastructure damage, and threaten important cultural heritage sites (Reimann et al., 2018; Vousdoukas et al., 2022). Storm surges arise from the interaction between the atmosphere and the sea. Essentially, the atmosphere exerts forces on the water body, causing sea levels to rise due to low atmospheric pressure systems and strong wind fields (Pirazzoli and Tomasin, 2022). The atmospheric pressure effect, known as the inverse barometer effect or static amplification, typically contributes 10 to 15% of the total storm surge magnitude (World  
40 Meteorological Organization, 2011). The second and more significant part of the storm surge, called dynamic amplification or wind setup, arises from tangential wind stress associated with the weather system's wind field acting on the ocean surface (Chaumillon et al., 2017).

Numerical simulations play a pivotal role in unraveling the complexities of physical phenomena, such as storm surges (Park  
45 et al., 2022). They offer invaluable insights into various processes and greatly contribute to building extensive databases for further analysis and comprehension. Concerning storm surge, this refers to a complex oceanographic phenomenon that demands accurate oceanic and atmospheric data for precise representation. Due to diverse orographic configurations, atmospheric models often exhibit significant errors, necessitating the utilization of local-scale models with high resolution (Umgiesser et al., 2021). Additionally, the intricate coastal and bathymetric features and interactions pose challenges for  
50 existing hydrodynamical models to fully capture the relevant dynamics, partly due to their low resolution (Mentaschi et al., 2015; Toomey et al., 2022).

On the other hand, the utilization of unstructured grid models enables a more accurate portrayal of coastal dynamics, considering the intricacies of bathymetry and shoreline configurations (Federico et al., 2017). This approach offers the  
55 advantage of employing higher resolution at the coastlines while maintaining more modest resolution in deeper waters (Ferrarin et al., 2019). Unstructured meshes offer flexibility in resolving basin geometry, allowing for local refinement of computational domains to simulate regional dynamics on a global mesh with coarse resolution. This flexibility is particularly valuable for coastal applications, where computational domains encompass complex coastlines and varying scales, ranging from basin size to details of river estuaries or riverbeds (Danilov, 2013). Over recent years, unstructured grid models have increasingly  
60 emerged as alternatives to regular grids for large-scale simulations (e.g. Mentaschi et al., 2020; Muis et al., 2016; Vousdoukas et al., 2018; Fernández-Montblanc et al., 2020; Saillour et al., 2021; Wang et al., 2022; Zhang et al., 2023; Mentaschi et al., 2023), with established circulation unstructured models like ADCIRC (Luettich et al., 1992; Pringle et al., 2021), the Finite-Volume Coastal Ocean Model (FVCOM, Chen et al., 2003), the Semi-implicit Cross-scale Hydroscience Integrated System Model (SCHISM, Zhang and Baptista, 2008; Zhang et al., 2016), the System of Hydrodynamic Finite Element Modules

65 (SHYFEM, Umgiesser et al., 2004; Bellafiore and Umgiesser, 2010; Micaletto et al., 2022), the model TELEMAC (Hervouet and Bates, 2000), Delft3D-FM (Deltares: Delft, 2024), among others.

**Commented [RC1]:** Reply Referee #1, Comment #1

In this study, we developed numerical simulations of storm surge in the Northern Adriatic Sea with two main objectives: first, to generate long-term databases of storm surge with a focus on accurately representing extreme values, and second, to analyze the ability of different metrics to capture the skill of the model. The Northern Adriatic Sea is a semi-enclosed body of water characterized by intricate bathymetry. The region's coastline exhibits distinct features, with the western coastline being relatively smooth and sandy, while the eastern coastline is fragmented and rocky, dotted with numerous islands. Both bathymetry and the configuration of the coastline significantly influence the physical processes occurring along the coast (Bellafiore and Umgiesser, 2010). The semi-enclosed nature of the Adriatic Sea predisposes it to experiencing intense storm surge events, leading to anomalous increases in sea level. These events are typically driven by local low-pressure system cyclogenesis and the associated strong winds, which are influenced by the region's orographic features (Umgiesser et al., 2021).

**Commented [RC2]:** Reply Referee #2, Comment #1

The application of numerical tools to study storm surge in the Northern Adriatic Sea has garnered significant attention over the years, primarily due to its status as a high-risk area with unique cultural and environmental heritage, as well as significant economic activities (Ferrarin et al., 2020). Previous efforts in this field have included predictive models projecting future storm scenarios (Yu et al., 1998), long-term numerical simulations (Lionello et al., 2010), analyses of storm events and use various atmospheric forcings (De Vries et al., 1995; Zampato et al., 2006; Medugorac et al., 2018), investigations into seiches influence and data assimilation impacts (Bajo et al., 2019), and storm surge ensemble prediction systems for lagoons (Alessandri et al., 2023).

85 In this study, the numerical simulations are based on a long-term ocean circulation downscaling carried out with the SHYFEM model, which is an unstructured-grid finite element hydrodynamic open-source code that solves the Navier-Stokes equations with hydrostatic and Boussinesq approximations (Umgiesser et al., 2004; Micaletto et al., 2022). The model has been already implemented in operational (Federico et al., 2017) and relocatable (Trotta et al., 2016) forecasting frameworks, and for storm surge events (Park et al., 2022; Alessandri et al., 2023). The choice of SHYFEM is driven by its flexibility in handling complex bathymetry and irregular coastlines through its unstructured-grid framework, allowing for higher resolution in critical areas. Additionally, its successful implementation in operational and relocatable forecasting frameworks, and storm surge events, confirms its reliability for this study. The simulations consider different setups to explore the influence of different atmospheric forcings and model configurations on the model's skill. Regarding model configurations, both barotropic and baroclinic simulations were conducted to compare potential differences between these two widely used approaches, as covered in the literature for the proper representation of storm surge (e.g. Weisberg and Zheng, 2008; Staneva et al., 2016; Hetzel et al., 2017; Ye et al., 2020; Muñoz et al., 2022). Furthermore, we focus on the use of different metrics and their ability to provide reliable indications of the model's performance, which is an essential aspect in assessing model skill and to select the best model

**Commented [RC3]:** Reply Referee #2, Comment #2

**Commented [RC4]:** Reply Referee #1, Comment #8

**Commented [RC5]:** Reply Referee #3, Comment #1

configuration. In addition to classical metrics such as the Pearson correlation coefficient and Root-Mean Squared Error (RMSE), two customized versions of the Mean Absolute Deviation (MAD) are introduced. These tailored metrics incorporate observed and simulated percentiles, ranging from 0 to 100%, to ensure accurate representation of extreme values during the performance evaluation.

The paper is organized as follows. Materials and methods are described in Section 2, including the description of the two atmospheric databases considered for the simulations, the model setup, and the procedures to carry out the performance evaluation. Section 3 shows the main results of the comparisons between observed and simulated storm surge. The paper continues with a discussion of the results on Section 4. Finally, the conclusion on Section 5 summarizes the key points of the study.

## 2 Materials and methods

### 2.1 Atmospheric forcing

In this study, we utilized two distinct atmospheric databases to force the circulation model, incorporating mean sea level pressure and wind fields. The first database is ERA5, the fifth generation of reanalysis data generated by the European Centre for Medium-Range Weather Forecasts (ECMWF). ERA5 builds upon the Integrated Forecasting System (IFS) Cy41r2, which became operational in 2016, providing hourly output with a horizontal resolution of  $0.25^\circ \times 0.25^\circ$  for atmospheric variables (Hersbach et al., 2020). ERA5 is relatively high resolution and accurate for a global reanalysis, although it is known to be affected by negative biases at high percentiles, particularly when compared with measured wind speeds (Pineau-Guillou et al., 2018; Vannucchi et al., 2021; Benetazzo et al., 2022; Gumuscu et al., 2023).

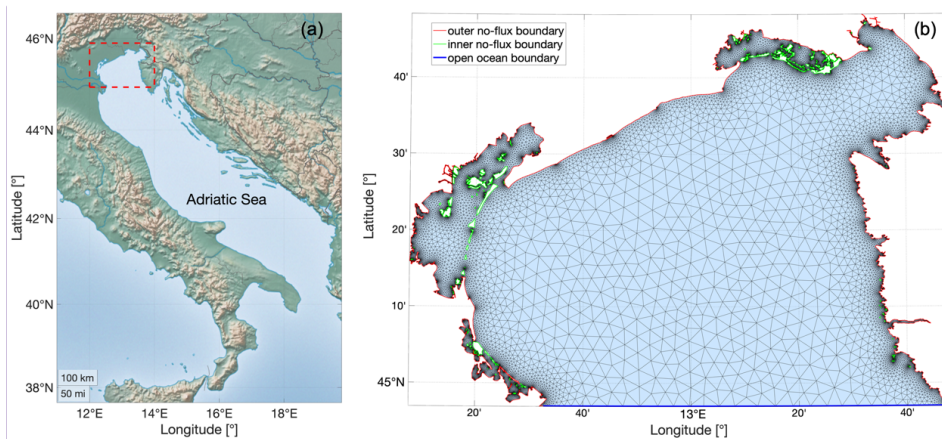
Since ERA5 is relatively coarse for local studies and exhibits significant underestimation of extremes, we employed an alternative approach using a high-resolution (3.3 km) atmospheric downscaling developed by the University of Genoa (UniGe). Wind forcing was derived from 10 m wind fields via the Weather Research and Forecast (WRF-ARW) model v3.8.1, allowing for improved representation of small-scale forcings and physics. The computational domains comprised a 10 km resolution grid covering the Mediterranean, Northern Africa, and Southern Europe (A10), and a 3.3 km grid over the Tyrrhenian Basin and Northern Adriatic basin (A3), nested within A10. Initial conditions were obtained from the Climate Forecast System Reanalysis (CFSR) data, known for reliability but occasionally underestimating extreme events (Saha et al., 2010). WRF simulations were conducted for 24 hours with hourly outputs, employing established physical parameterization schemes to ensure accuracy across various atmospheric conditions. For further details, readers are referred to Mentaschi et al. (2015).

Commented [RC6]: Reply Referee #1, Comment #10

Commented [RC7]: Reply Referee #2, Comment #3

## 2.2 Model setup

The SHYFEM model utilizes staggered finite elements in an unstructured Arakawa B horizontal grid, with the vertices of the triangle elements referred to as nodes. Vectors (velocity) are calculated at the center of each element, while scalars (temperature, salinity, and water levels) are determined at nodes (Federico et al., 2017). The unstructured grid for the simulations in this study was generated using the OceanMesh2D tool (Roberts et al., 2019) with a horizontal resolution of 3 km on the open ocean boundary and 50 m in the coastline (Fig. 1.a). The General Bathymetric Chart of the Oceans (GEBCO) dataset (Weatherall et al., 2015) was used, incorporating a high-resolution coastline from the European Environmental Agency. However, due to identified overestimations in water depth in the Venice and Marano lagoons from GEBCO bathymetry, adjustments were made based on the contributions from Fagherazzi et al. (2007), Lovato et al. (2010), Zaggia et al. (2017) for the Venice lagoon and Petti et al. (2019) and Bosa et al. (2021) for the Marano lagoon.



140 **Figure 1: (a) Location of study area, marked with dashed red line; (b) Unstructured grid for study area, in which the blue line**  
145 **represents the location of the open boundary condition, the red line the coastline, and the green lines the coastline formed by**  
**islands.**

As initial and open ocean boundary conditions, sea level residuals, current velocity, temperature, and salinity from the Copernicus Mediterranean Sea Physics reanalysis (Escudier et al., 2021) were considered. Tides with hourly resolution from the Finite Element Solution (FES) 2014 (Lyard et al., 2021) were also included to account for the total sea level in the simulations. Specifically, the constituents included for the tide reconstruction are SA, SSA, O1, P1, S1, K1, N2, M2, MKS2, S2, R2, K2, M3, M4, and MS4, which were selected based on preliminary harmonic analysis applied to sea level observation data in the locations specified in Section 2.2.

Commented [RC8]: Reply Referee #2, Comment #4

Commented [RC9]: Reply Referee #1, Comment #11

Commented [RC10]: Reply Referee #1, Comment #2

150 Two model configurations were considered: a) barotropic (BT) and b) baroclinic, employing 33 vertical levels with layer thickness of 1 m up to 10 m depth and then 2 m up to a maximum depth of 60 m (BC). To determine vertical viscosities and diffusivities, we utilize a k-ε turbulence scheme derived from the General Ocean Turbulence Model (GOTM) model (Burchard and Petersen, 1999). For wind stress at the air-sea interface a constant wind drag coefficient of  $2.5 * 10^{-3}$  was employed, following the works from Orlic' et al. (1994) and Zampato et al. (2007). The bottom stress is determined through the quadratic  
 155 formulation:

$$\tau_{xz}^{zN} = \frac{C_B}{H_N^2} |\overline{U_N}| U_N \quad \tau_{yz}^{zN} = \frac{C_B}{H_N^2} |\overline{V_N}| V_N \quad (1)$$

Where  $\tau_{xz}^{zN}$  and  $\tau_{yz}^{zN}$  are the turbulent shear stresses at the bottom interface of the deepest layer,  $H_N$  is bottom layer thickness,  
 160  $U_N$  and  $V_N$  the zonal and meridional transports of the bottom layer.  $C_B$  is the bottom drag coefficient defined as:

$$C_B = \left( \frac{0.4}{\ln\left(\frac{\lambda_B + 0.5H_N}{\lambda_B}\right)} \right)^2 \quad (2)$$

Where  $\lambda_B$  is the bottom roughness length expressed in m, which in this study remains constant at 0.01 m. For further details,  
 165 readers are referred to Maicu et al. (2015).

Commented [RC11]: Reply Referee #2, Comment #5

The simulation period extends from 1987 to 2020, with hourly output. Three combinations of atmospheric forcing and configuration are considered here: 1) barotropic forced by ERA5 (BT-ERA5), 2) baroclinic forced by ERA5 (BC-ERA5), and 3) baroclinic forced by UniGe (BC-UniGe).

### 170 2.3 Model performance evaluation

The model output was compared with observations from tide gauges located in the Northern Adriatic Sea. The observational data were acquired from the Italian National Institute for Environmental Protection and Research (ISPRA), the Civil Protection of the Friuli-Venezia Giulia Region, and Raicich (2023). Table 1 summarizes the locations considered, and the available time spans for comparison that match with the simulation timespan. Fig. 2 shows the locations considered for comparison between  
 175 measured and simulated storm surge, together with the bathymetry used for the simulations.

Commented [RC12]: Reply Referee #1, Comment #12

Location	Lon [°]	Lat [°]	Start date	End date
ISMAR-CNR research platform "Aqua Alta" (hereafter CNR platform)	12.53	45.31	01-01-1987	31-12-2020
Punta della Salute	12.33	45.43	01-01-1987	31-12-2020
Caorle	12.86	45.59	01-01-2000	31-12-2020
Grado	13.38	45.68	01-01-1991	31-12-2020
Monfalcone	13.54	45.78	01-01-2008	31-12-2020
Trieste	13.76	45.64	01-01-1987	31-12-2020

Commented [RC13]: Minor changes on table format

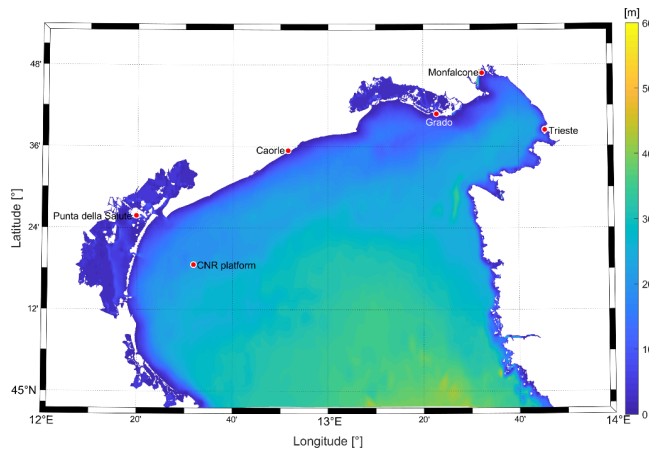


Figure 2: Tide gauges locations and bathymetry (depth values on positive).

185 Both the model output and the observations were processed as follows to enable their intercomparability. To start, both measurement and simulation were centered with a zero mean and then detrended. This approach mitigates possible effects of unmodulated land motion (Chepurin et al., 2014) and ensures that extreme values across the years can be considered as homogeneous and can be compared despite relative sea level changes (Ferrarin et al., 2022). Harmonic analysis was performed for each calendar year on the detrended sea levels using the T-Tide MATLAB package (Pawlowicz et al., 2022), and the non-tidal residual was obtained as the arithmetic difference between sea level and tides (Tiggeloven et al., 2021). Performing yearly harmonic analysis reduces timing errors that could cause tidal energy to seep into the non-tidal residual (Merrifield et al., 2013).

Commented [RC14]: Reply Referee #1, Comment #13

195 Finally, to obtain the pure storm surge (hereafter also called “surges”), a low-pass filter is applied to the non-tidal residual, following the work from Park et al. (2022). In this study, we consider a cut-off period of 13 hours for the filter based on the mixed-semidiurnal tidal regime around the Northern Adriatic Sea (Lionello et al., 2021).

The performance evaluation of the simulations relies on the computation of statistical metrics of hourly data, which encompass the entire dataset, as well as values exceeding the 99th percentile from the cumulative distribution of measured data at each location. The following metrics are considered:

Pearson correlation:

$$\rho = \frac{1}{N-1} \sum_{i=1}^N \left( \frac{S_i - \mu_S}{\sigma_S} \right) \left( \frac{O_i - \mu_O}{\sigma_O} \right) \quad (3)$$

205

Where  $S_i$  and  $O_i$  are the  $i$ th simulated and observed data respectively,  $N$  is the sample size,  $\mu$  and  $\sigma$  are the mean and standard deviations of  $S$  and  $O$ . A value closer to one identifies a better performance.

Root-Mean Squared Error (RMSE):

210

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (S_i - O_i)^2} \quad (4)$$

A value closer to zero indicates a better performance.

215 Bias:

$$Bias = S - O \quad (5)$$

220 Where  $S$  and  $O$  are the average simulation and observation values respectively. A value closer to zero identifies a better performance, negative values indicate underestimation, and positive values indicate overestimation from the simulations. Given that both observed and simulated data were detrended and had their mean removed, bias was solely applied to the analysis of values exceeding the 99th percentile.

Slope of linear fit between observations and simulation:

225



$$S = m O + b \quad (6)$$

Where the slope is given by the coefficient  $m$ . A value closer to one indicates a better performance.

230 Mean Absolute Deviation (MAD):

$$MAD = |S - O| \quad (7)$$

A value closer to one indicates a better performance.

235

Additionally, with the aim of considering the representation of extremes by the simulations, we introduce two new metrics based on customized versions of the Mean Absolute Deviation:

MAD of the percentiles (MADp):

240

$$MADp = |S_{prc} - O_{prc}| \quad (8)$$

Where  $S_{prc}$  and  $O_{prc}$  are the simulation and observation percentile values, considered from 0 to 100%, every 1%. The MADp metric provides a comprehensive assessment of simulation model performance by comparing percentile values derived from simulations ( $S_{prc}$ ) with those observed ( $O_{prc}$ ). This evaluation encompasses the entire distribution, from the lowest to the highest percentiles, allowing to gauge the model's accuracy across a range of scenarios. MADp is particularly valuable for its sensitivity to systematic errors, such as persistent underestimation of high percentiles, which can significantly impact the reliability of simulation results. By penalizing these systematic errors, MADp highlights areas where improvements in the simulation model are necessary to better align with observed data. Lower MADp values indicates closer agreement between simulations and observations.

250

Corrected MAD (MADc):

$$MADc = |S - O| + MADp \quad (9)$$

255

In this indicator we exploit the ability of the “traditional” MAD to capture the model’s skill but reduce its strong penalization of the phase error or timing error (i.e. the reproduction by the model of peaks shifted in space-time) by adding the MAD (MADp) on the percentiles previously defined. MAD measures the average absolute difference between simulated and

**Commented [RC15]:** Reply Referee #1, Comment #3a.

Reply Referee #2, Comment #6

observed values, while MADp evaluates the average percentage deviation between them. By combining these two components, MADc provides a comprehensive evaluation of the simulation model's performance, considering both the magnitude and percentage deviations. A lower MADc value indicates better agreement between simulated and observed values, reflecting higher accuracy and reliability of the simulation model.

To quantify phase errors between observations and simulations, peaks in the hourly time series were identified using Matlab's 'find peaks' function for both observed and simulated data. The phase error was then calculated by measuring the time difference, in hours, between the occurrence of each peak in the observations and the corresponding peak in the simulations. This approach provided a direct assessment of the model's accuracy in capturing the timing of key events, such as storm surges.

The proposed metrics were also validated using an idealized time series. A sinusoidal time series was generated to represent an observed parameter. Two simulated time series were then created: one with the same amplitude as the observation but shifted in time (introducing a phase error), and the other with the same phase as the observation but with half the amplitude. Various metrics were calculated and plotted on scatter plots (Fig. 1.S). The results indicated better performance for the simulation that underestimated the observations when assessed with Pearson correlation, RMSE, and MAD. In contrast, the time series that accurately captured the amplitude was penalized for the phase error, which negatively affected its performance on these metrics. However, the proposed MADp and MADc metrics identified it as the better model.

### 3 Results

The Probability Distribution Estimates (PDE) and Empirical Cumulative Distribution Functions (ECDF), available in Fig. 2S to 7S, show that BC-UniGe better represents the higher values of storm surge when compared with observations, particularly when considering values above the 99th percentile. However, some overestimations are noticeable in Caorle and Monfalcone with BC-UniGe. In contrast, simulations with ERA5 forcing tend to underestimate these higher values, which is more noticeable for BT-ERA5.

The performance evaluation shows that, if the model performance is assessed in terms of Pearson correlation, RMSE, and MAD, the surges simulated with the ERA5 forcing fit better to the measured data (Fig. 3). The Pearson correlation coefficients obtained range between 0.8 and 0.9 in all locations for the three simulations, with maximum of 0.842 with BT-ERA5 in Grado (Fig. 3.d). Regarding the RMSE, mean values of 0.077 m for BT-ERA5, 0.075 m for BC-ERA5, and 0.079 m for BC-UniGe were obtained, with a minimum of 0.072 m (BT-ERA5 in Grado, Fig. 3.d) and a maximum of 0.094 m (BC-UniGe in Monfalcone, Fig. 3.e). Similar results are obtained for MAD, which shows better performance for the simulations with ERA5 forcing at all locations. Only in Trieste does BC-UniGe achieve the same performance as BC-ERA5 for this metric. Despite the aforementioned, the best performance is achieved by BC-UniGe in the linear fit slope, with values above 0.8 in all locations

**Commented [RC16]:** Reply Referee #1, Comment #6

**Commented [RC17]:** Reply Referee #2, Comment #6

**Commented [RC18]:** Reply Referee #2, Comment #7

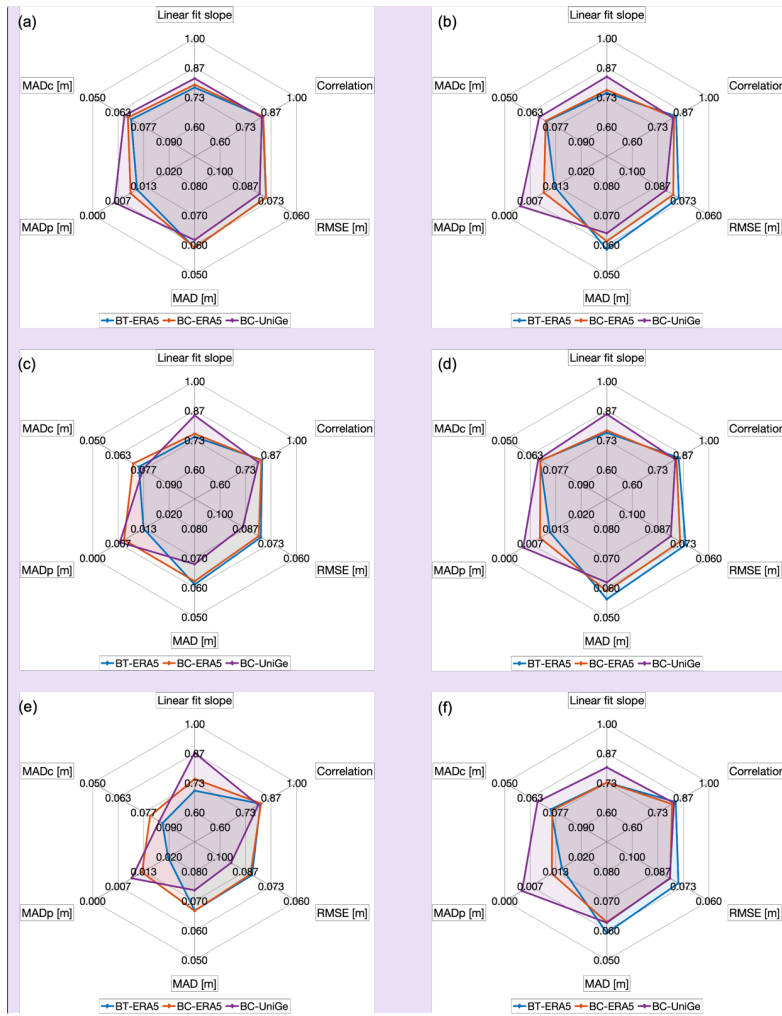
**Commented [RC19]:** Reply Referee #1, Comment #3a  
Reply Referee #2, Comment #6

**Commented [RC20]:** Reply Referee #1, Comment #3a  
Reply Referee #2, Comment #6

and a maximum of 0.869 in Monfalcone (Fig. 3.e). For this parameter, the less favorable performance is obtained with BT-ERA5 in all locations.

295 For MADp, the best performance is achieved by BC-UniGe in all locations, with a mean value of 0.004 m, while less favorable results are obtained with BT-ERA5, with a mean of 0.011 m. Similar results were obtained for MADc, except in Caorle (Fig. 3.c) and Monfalcone (Fig. 3.e), where BC-ERA5 showed better performance, likely due to overestimation in the mentioned sites. These results underscore the importance of considering percentiles as part of the performance evaluation. BC-UniGe simulations demonstrate an improvement in representing extreme values, showing a better fit of the highest percentiles, which can be noticed in Fig. 4 and Fig. 5. Additionally, these figures indicate that BC-UniGe simulations produce greater dispersion  
300 of data, likely due to a more frequent occurrence of phase error, which was quantified as 3.1% higher than in BT-ERA5, and 4.5% higher than in BC-ERA5. However, they also exhibit a better fit of the linear regression and a more accurate representation of extreme values compared to BC-ERA5, which fail to represent the most extreme events in each location.

**Commented [RC21]:** Reply Referee #1, Comment #6



305 **Figure 3:** Radar charts of evaluation metrics for the total amount of data in all locations. (a) CNR platform; (b) Punta della Salute; (c) Caorle; (d) Grado; (e) Monfalcone; (f) Trieste. For RMSE, MADp and MADc a reverse axis is used, this ensures that simulations covering a larger area on each metric represent a better performance (i.e. values on the fringe refer to better performance).

**Commented [RC22]:** Reply Referee #1, Comment #3a

Reply Referee #2, Comment #6

**Commented [RC23]:** Reply Referee #2, Comment #8

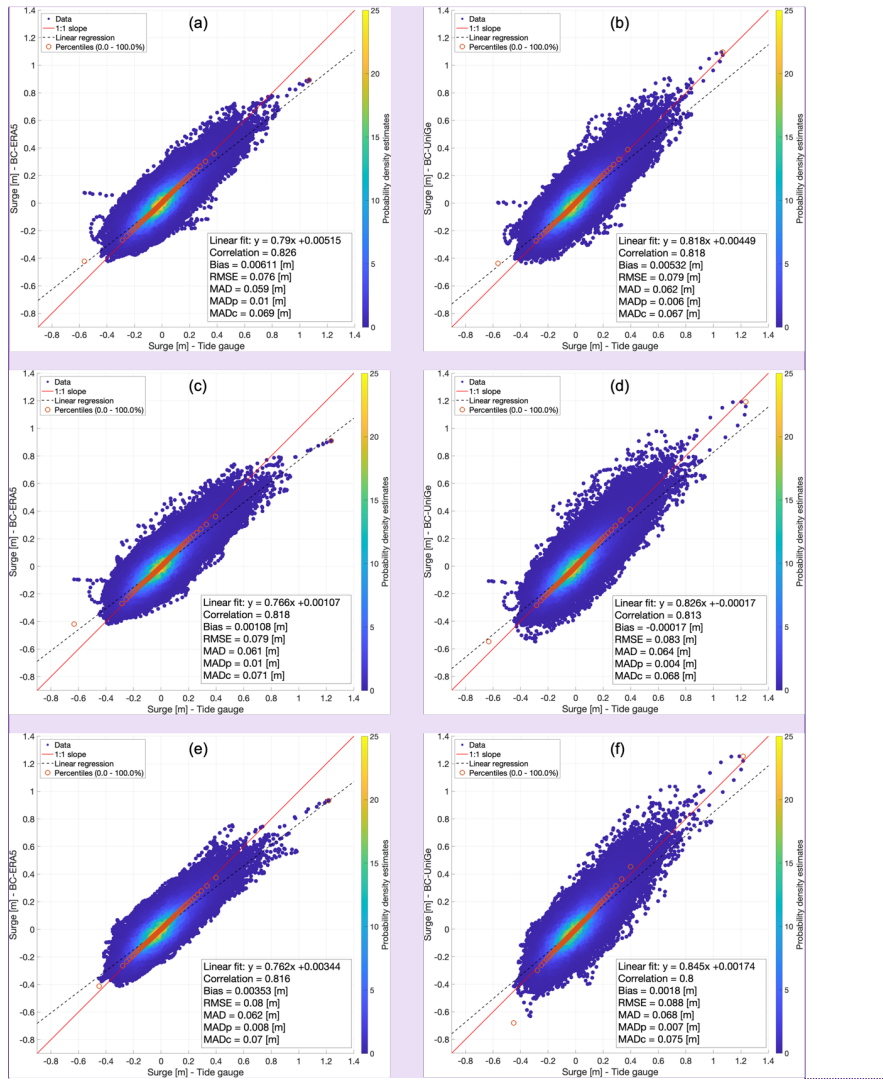


Figure 4: Scatter plots between tide gauges and baroclinic simulations. CNR platform: BC-ERA5 (a), BC-UniGe (b); Punta della Salute: BC-ERA5 (c), BC-UniGe (d); Caorle: BC-ERA5 (e), BC-UniGe (f).

Commented [RC24]: Reply Referee #1, Comment #3a

Reply Referee #2, Comment #6

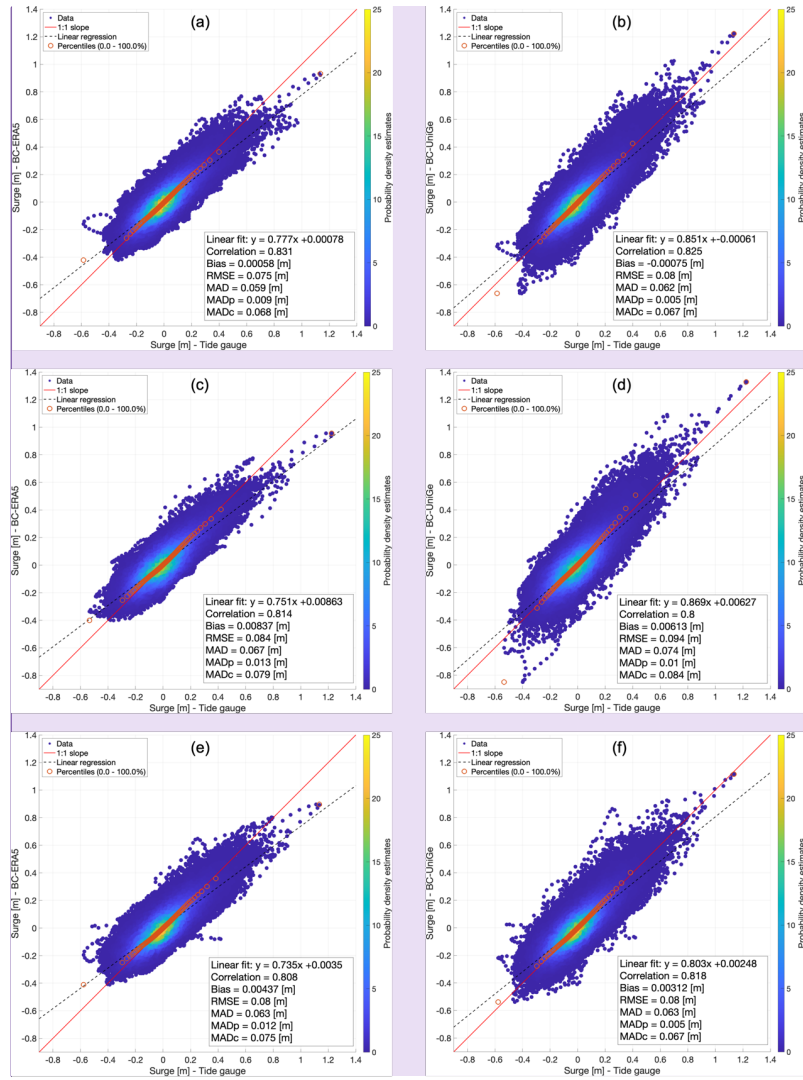


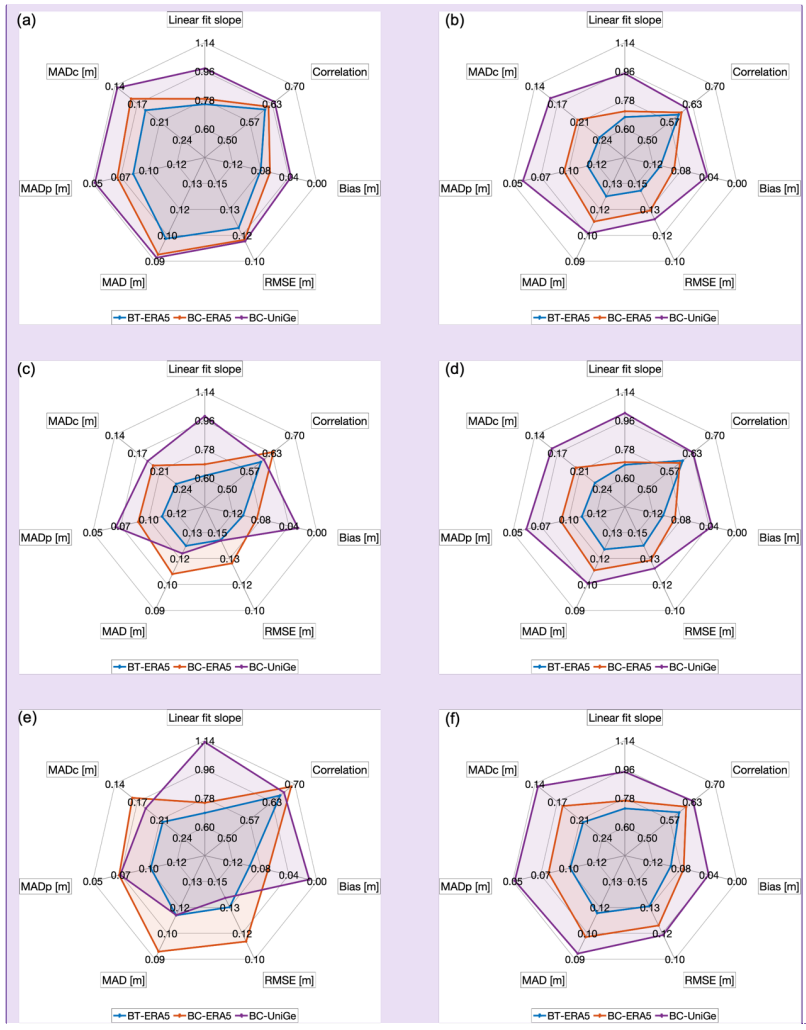
Figure 5: Scatter plots between tide gauges and baroclinic simulations. Grado: BC-ERAS (a), BC-UniGe (b); Monfalcone: BC-ERAS (c), BC-UniGe (d); Trieste: BC-ERAS (e), BC-UniGe (f).

315 The results of the error metrics for surge values above the 99th percentile, represented using radar charts (Fig. 6), confirm that, in general, better performance is observed with BC-UniGe, while less favorable results are obtained for BT-ERA5. Although the transition from barotropic to baroclinic configuration indicates an improvement in the representation of extremes (Weisberg and Zheng, 2008; Staneva et al., 2016; Hetzel et al., 2017; Ye et al., 2020; Muñoz et al., 2022), the utilization of UniGe forcing represents the best improvement across practically all metrics. Only in Caorle (Fig. 6.C) and Monfalcone (Fig. 6.E) does BC-  
320 ERA5 show better Pearson correlation, RMSE, and MAD; additionally, in the latter, MADc exhibits better performance for that simulation, likely due to overestimation of the peaks by BC-UniGe in Monfalcone. In the other locations, it's evident that BC-UniGe performs better in representing the highest storm surge values.

In order to show the capacity of the different model configurations to represent certain known storm events at each location, Fig. 7 shows time series of different storm surge events at each location. These extreme events were chosen according to the contributions of Lionello et al. (2012), Medugorac et al. (2018), Ferrarin et al. (2020), Umgiesser et al. (2021), and Giesen et al. (2021). As mentioned before, the incorporation of the UniGe forcing implies a significant improvement in the representation of extreme events, clearly evident in the peak values of the storm surge. Despite, an overestimation of some surge peaks is also observed in the events chosen at Punta della Salute (Fig. 7.B), Caorle (Fig. 7.C), and Monfalcone (Fig. 7.E) with BC-  
330 UniGe. On the other hand, a systematic underestimation of extremes obtained in simulations with ERA5 forcing is notable on every surge peak.

**Commented [RC26]:** Reply Referee #1, Comment #3a

Reply Referee #2, Comment #6



335

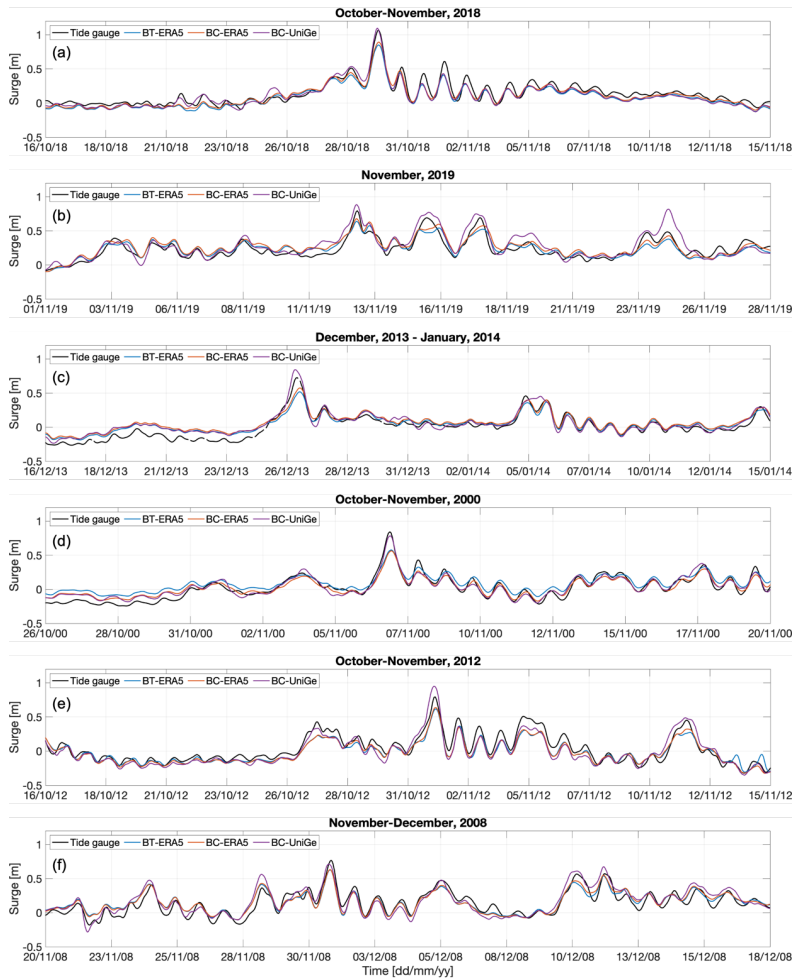
Figure 6: Radar charts of evaluation metrics for surge values above the 99th percentile of the cumulative distribution at each location. (a) CNR platform; (b) Punta della Salute; (c) Caorle; (d) Grado; (e) Monfalcone; (f) Trieste. Bias is represented by absolute value. Also, for RMSE, Bias, and MADp and MADc a reverse axis is used, this ensures that simulations covering a larger area on each metric represent a better performance (i.e. values on the fringe refer to better performance).

Commented [RC27]: Reply Referee #1, Comment #3a

Reply Referee #2, Comment #6

Commented [RC28]: Reply Referee #2, Comment #8





340 **Figure 7: Time series of different storm surge events in all the locations, tidal gauge versus model. (a) CNR platform; (b) Punta della Salute; (c) Caorle; (d) Grado; (e) Monfalcone; (f) Trieste.**

#### 4 Discussion

345 The utilization of different atmospheric forcing databases has revealed significant implications for the representation of storm surge in numerical simulations. Given the direct influence of wind speed and sea level pressure on this phenomenon, as represented in both forcings databases, the resulting model performances present significant differences. While simulations using ERA5 forcing generally show slightly better performance on traditional metrics such as RMSE, MAD, and Pearson correlation coefficient, a more detailed analysis reveals that using the UniGe forcing results in better performance, especially the extreme values when considering additional metrics.

350 Simulations using ERA5 forcing tend to underestimate the highest surge values, primarily due to a corresponding underestimation of extreme wind speed by this database, a variable crucially linked to surge amplitude (Campos et al., 2022). Despite this, metrics such as Pearson correlation, RMSE, and MAD generally indicate better performance for ERA5 simulations. Conversely, the utilization of UniGe forcing shows an improvement in representing the peaks of storm surge events (with the noticeable exception of Monfalcone, where the extremes are overestimated, and where MADp present similar values for BC-ERA5 and BC-UniGe). These results demonstrate that the increase in atmospheric forcing resolution does not consistently translate into better values of all the statistical metrics.

360 It is important to recognize that identifying the optimal model configuration cannot rely solely on a few statistical metrics. As outlined in section 3 no single simulation emerges as superior across all metrics and locations. While ERA5 simulations may demonstrate better performance on RMSE, Pearson correlation, and MAD, BC-UniGe exhibits superior performance in terms of the slope of the linear fit, MADp and MADc.

365 From an epistemic point of view BC-UniGe is a significantly more sophisticated model compared to BT-ERA5. Not only does it employ a higher resolution forcing. It also takes into account the baroclinicity and the vertical motion within the water column, whereas the barotropic configuration of BT-ERA5 approximates the ocean as a 2D sheet only subject to vertically uniform motions and waves. This suggests that widespread indicators such as RMSE, Pearson correlation, and MAD, which in this case identify BT-ERA5 as the best model, should not be considered as the sole source of information in model skill assessment, since a higher resolution forcing and a baroclinic setup are known in literature to better capture the variability of the sea levels (Weisberg and Zheng, 2008; Hetzel et al., 2017; Muñoz et al., 2022).

375 Similar results were found by Zampato et al. (2006) using SHYFEM with three different forcings for wind and atmospheric pressure fields: ECMWF global model, high-resolution LAMI model and satellite QuickSCAT. In this work, the authors found well correlated sea levels with observations near Venice using the ECMWF forcings, but underestimation on highest values. On the other hand, simulations driven by the high-resolution model (LAMI) succeeded in simulating the storm surge, giving

**Commented [RC29]:** Reply Referee #1, Comment #3a

Reply Referee #2, Comment #6

**Commented [RC30]:** Reply Referee #1, Comment #3a

Reply Referee #2, Comment #6

a good reproduction of the sea level peaks. Nevertheless, the correlation with observed data was lower than in the case of ECMWF forcing.

380 The complexity in simulations performance evaluation is echoed in the work of Mentaschi et al. (2013), who caution against over-reliance on metrics like RMSE, NRMSE (Normalized RMSE), and SI (Scatter Index) as indicators of model performance. These metrics may not fully capture the intricacies of natural processes such as atmospheric dynamics, ocean circulation, or wave generation and propagation. These authors mention that the RMSE and its variations tend to assume typical values of the best performance for simulations that underestimate the physical process of interest. The discrepancy between metrics and  
385 the representation of extremes highlights the need for a comprehensive understanding of model performance beyond traditional statistical measures.

This results on performance evaluation are usually related to phase error in high-resolution models and RMSE “double penalty”. The phase error refers to a discrepancy between the timing or phase of a simulated event and its actual occurrence  
390 on measured data. In the context of atmospheric models, phase errors can manifest as delays or advances in the timing of weather events, such as the onset of precipitation, the movement of storm systems, or the arrival of fronts. Double penalty refers to a situation where the errors in the model output are penalized twice, in indicators such as RMSE and MAD, once for missing the observations and again for giving a false alarm (e.g. Gilleland et al., 2009). This is a well-known problem during performance evaluation of numerical models and different contributions have sought to overcome it, with approaches  
395 specialized in atmospheric and oceanographic fields (e.g., Ebert and McBride, 2000; Zingerle and Nurmi, 2008; Roberts and Lean, 2008; Mittermaier, 2014; Skok and Roberts, 2016; Crocker et al., 2020).

In RMSE, “double penalty” is further amplified compared to MAD, as the penalizations due to the peak mismatch are squared. This means that phase errors have a disproportionately large impact on RMSE. A more sophisticated model may be better able  
400 to capture the magnitude of the peaks, but as it is more prone to phase error compared to low-resolution ones this ability will be doubly penalized. This is the reason why a less sophisticated model employing a low-resolution forcing (BT-ERA5) appears to out-perform the other two in terms of RMSE. Conversely, MAD, although it also experiences a form of “double penalty,” reduces the impact of this effect compared to RMSE. As a result, the performance differences between simulations, particularly above the 99th percentile, are generally more pronounced for MAD than for RMSE, better highlighting the superiority of BC-  
405 UniGe. This enhanced differentiation is likely due to MAD's linear weighting of errors, which reduces the inflated impact of large deviations that characterize RMSE.

In other words, RMSE tends to be better for “blurring” models, whereas high-resolution models, known to be more capable of reproducing small-scale dynamics (e.g. BC-UniGe), perform worse in terms of RMSE due to phase error (Crocker et al., 2020).

**Commented [RC31]:** Reply Referee #1, Comment #3a

Reply Referee #2, Comment #6

410 Although in many aspects, capturing a peak with a phase error is preferable to missing the peak entirely, this does not lead to a reduction in the RMSE.

This limitation of RMSE also impacts the Pearson correlation. Indeed, RMSE can be decomposed into a bias component and a scatter component that depends solely on the Pearson correlation (Mentaschi et al., 2013, equation 8). All these considerations call for caution when claiming that one model outperforms another one just based on a better value of RMSE or MAD or Pearson correlation.

The MADc indicator was introduced here as a possible way to correct MAD to make it less prone to the double penalty effect. The incorporation in MADc of a term that takes into account the distribution of the data (the MAD of the percentiles MADp) rewards the ability of a high-resolution and more sophisticated model to reproduce the variability in the observations without systematic errors. In other words, MADc remains more resilient to phase errors compared to other metrics, ensuring that discrepancies in the timing of events do not unduly influence the assessment of model performance. The differences between the simulation metrics are generally in the range of millimeters when considering the overall data, but these differences are significant in relative terms. For the MADc metric, BC-UniGe shows improvements ranging from 1.3% (Grado) to 9.3% (Trieste) compared to BT-ERA5, and from 1.6% (Grado) to 10.3% (Trieste) compared to BC-ERA5. The improvements are even more notable when focusing on values above the 99th percentile, where BC-UniGe outperforms BT-ERA5 by 12% (Monfalcone) to 31.6% (Trieste), and BC-ERA5 by 4.1% (Caorle) to 20.2% (Trieste).

As shown in Section 3, some discrepancies were observed in Caorle and Monfalcone, where BC-ERA5 achieved better performance in terms of MADc. A possible explanation for this could be related to the location of the tide gauges at these sites. The tide gauge at Caorle is situated in a protected area inside the Livenza River, a location not fully represented by the simulations due to the resolution of the coastline, even though a high-resolution model data was used. A similar issue is found in Monfalcone, where the tide gauge is located in front of a breakwater not fully represented by the coastline used in the model. These factors could affect the signals obtained from observations and simulations, primarily due to local effects at the tide gauge locations.

## 5 Conclusions

In this study we developed high-resolution simulations of storm surge in the Northern Adriatic Sea spanning from 1987 to 2020, using the model SHYFEM, employing different forcing data and physical configurations. The comparative analysis of the results highlights nuanced differences in performance metrics, particularly concerning the representation of the extreme values. Traditional metrics like Pearson correlation, RMSE, and MAD favor a simulation (BT-ERA5) forced by a coarser database and employing a less sophisticated setup (barotropic). However, a closer examination and the use of different metrics

Commented [RC32]: Reply Referee #1, Comment #4a

Commented [RC33]: Reply Referee #1, Comment #4b

Commented [RC34]: Reply Referee #1, Comment #3a  
Reply Referee #2, Comment #6

tell a different story and allow to identify a baroclinic model forced by a high-resolution dataset (BC-UniGe) as better able to capture the variability of the water levels and, in particular, the extremes. This is because BC-UniGe is more prone to phase error than BT-ERA5, and is thus doubly penalized in indicators such as RMSE, MAD and Pearson correlation.

445

The corrected MAD (MADc) introduced in this study comes as a possible way to alleviate the double penalty, by adding a term that rewards the ability of a model to capture the distribution of the observations irrespective of the position of the peaks. In this study MADc is successful in identifying BC-UniGe as the best simulation in most locations. Even though this study has focused on the performance evaluation of storm surge, the analysis and proposed customized metrics (MADc and MADp) can be applied to any problem of validating a numerical model with observations by time-series comparison.

450

These findings suggest that simply having a lower RMSE is insufficient evidence to claim that one model is superior to another. RMSE, MAD and Pearson correlation are valuable indicators but should be used considering their limitations, and complemented by other metrics, qualitative assessment, and expert judgment.

#### 455 **6 Author contribution**

RC carried out the numerical simulations, post processing, performance evaluation of the simulations and prepared the manuscript. LM guided during numerical simulations, post processing, performance evaluation and contributed to the preparation of the manuscript. JA guided and supported numerical simulations. PC contributed to the performance evaluation and with the preparation of the manuscript. AM, FF, IF, and MV contributed during the preparation of the manuscript. MT contributed to the performance evaluation.

460

#### **7 Competing interests**

Co-author MT is employed by the company HS Marine SrL.

The remain authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

465

#### **References**

Alessandri, J., Pinardi, N., Federico, I., and Valentini, A.: Storm Surge Ensemble Prediction System for Lagoons and Transitional Environments. American Meteorological Society, 38. <https://doi.org/10.1175/WAF-D-23>, 2023.

**Commented [RC35]:** Reply Referee #1, Comment #3a

Reply Referee #2, Comment #6

**Commented [RC36]:** The list of references was edited based on the journal's format

- 470 Bajo, M., Medugorac, I., Umgiesser, G., and Orlić, M.: Storm surge and seiche modelling in the Adriatic Sea and the impact of data assimilation. *Quarterly Journal of the Royal Meteorological Society*. <https://doi.org/10.1002/qj.3544>, 2019.
- Bellafiore, D., and Umgiesser, G.: Hydrodynamic coastal processes in the north Adriatic investigated with a 3D finite element model. *Ocean Dynam.*, 60, 225–273. <https://doi.org/10.1007/s10236-009-0254-x>, 2010.
- 475 Benetazzo, A., Davison, S., Barbariol, F., Mercogliano, P., Favaretto, C., and Sclavo, M.: Correction of ERA5 Wind for Regional Climate Projections of Sea Waves. *Water (Switzerland)*, 14(10). <https://doi.org/10.3390/w14101590>, 2022.
- Bosa, S., Petti, M., and Pascolo, S.: Improvement in the sediment management of a lagoon harbor: The case of Marano  
480 Laganare, Italy. *Water*, 13. <https://doi.org/10.3390/w13213074>, 2021.
- Burchard, H., and Petersen, O.: Models of turbulence in the marine environment—a comparative study of two-equation turbulence models. *Journal of Marine Systems*, 21, 29–53, 1999.
- 485 Campos, R. M., Gramscianinov, C. B., de Camargo, R., and da Silva Dias, P. L.: Assessment and Calibration of ERA5 Severe Winds in the Atlantic Ocean Using Satellite Data. *Remote Sensing*, 14(19). <https://doi.org/10.3390/rs14194918>, 2022.
- Chaumillon, E., Bertin, X., Fortunato, A. B., Bajo, M., Schneider, J. L., Dezileau, L., Walsh, J. P., Michelot, A., Chauveau, E., Créach, A., Hénaff, A., Sauzeau, T., Waeles, B., Gervais, B., Jan, G., Baumann, J., Breilh, J. F., and Pedreros, R.: Storm-induced marine flooding: Lessons from a multidisciplinary approach. *Earth-Science Reviews*, 165, 151–184.  
490 <https://doi.org/10.1016/j.earscirev.2016.12.005>, 2017.
- Chen, C., Liu, H., and Beardsley, R. C.: An Unstructured Grid, Finite-Volume, Three-Dimensional, Primitive Equations Ocean Model: Application to Coastal Ocean and Estuaries. *Ocean. Technol.*, 20, 159–186, 2003.
- 495 Chepurin, G.A., Carton, J.A., and Leuliette, E.: Sea level in ocean reanalyses and tide gauges. *J. Geophys. Res. Oceans*, 119, 147–155. <https://doi.org/10.102/2013JC009365>, 2014.
- Crocker, R., Maksymczuk, J., Mittermaier, M., Tonani, M., and Pequignet, C.: An approach to the verification of high-  
500 resolution ocean models using spatial methods. *Ocean Science*, 16, 831–845. <https://doi.org/10.5194/os-16-831-2020>
- Danilov, S. (2013). Ocean modeling on unstructured meshes. *Ocean Modelling*, 69, 195–210. <https://doi.org/10.1016/j.ocemod.2013.05.005>, 2020.

- De Vries, H., Breton, M., De Mulder, T., Krestenitis, Y., Ozer, J., Proctor, R., Ruddick, K., Salomon, J. C., and Voorrips, A.:  
505 A comparison of 2D storm surge models applied to three shallow European seas. *Environmental Software*, 10, 23–42, 1995.
- Deltares: Delft: Delft3D-FLOW User Manual, 2024.
- Ebert, E. E., and McBride, J. L.: Verification of precipitation in weather systems: determination of systematic errors. *Journal*  
510 *of Hydrology*, 239, 179–202, 2000.
- Escudier, R., Clementi, E., Cipollone, A., Pistoia, J., Drudi, M., Grandi, A., Lyubartsev, V., Lecci, R., Aydogdu, A., Delrosso,  
D., Omar, M., Masina, S., Coppini, G., and Pinardi, N.: A High Resolution Reanalysis for the Mediterranean Sea. *Frontiers in*  
*Earth Science*, 9. <https://doi.org/10.3389/feart.2021.702285>, 2021.
- 515 Fagherazzi, S., Palermo, C., Rulli, M. C., Carniello, L., and Defina, A.: Wind waves in shallow microtidal basins and the  
dynamic equilibrium of tidal flats. *Journal of Geophysical Research: Earth Surface*, 112(2).  
<https://doi.org/10.1029/2006JF000572>, 2007.
- 520 Federico, I., Pinardi, N., Coppini, G., Oddo, P., Lecci, R., and Mossa, M.: Coastal ocean forecasting with an unstructured grid  
model in the southern Adriatic and northern Ionian seas. *Natural Hazards and Earth System Sciences*, 17(1), 45–59.  
<https://doi.org/10.5194/nhess-17-45-2017>, 2017.
- Fernández-Montblanc, T., Vousdoukas, M. I., Mentaschi, L., and Ciavola, P.: A Pan-European high resolution storm surge  
525 hindcast. *Environ. Int.*, 135. <https://doi.org/10.1016/j.envint.2019.105367>, 2020.
- Ferrarin, C., Davolio, S., Bellafiore, D., Ghezzi, M., Maicu, F., Mc Kiver, W., Drofa, O., Umgieser, G., Bajo, M., De Pascalis,  
F., Malguzzi, P., Zaggia, L., Lorenzetti, G., and Manfè, G.: Cross-scale operational oceanography in the Adriatic Sea. *Journal*  
*of Operational Oceanography*, 12, 86–103. <https://doi.org/10.1080/1755876X.2019.1576275>, 2019.
- 530 Ferrarin, C., Lionello, P., Orlić, M., Raicich, F., and Salvadori, G.: Venice as a paradigm of coastal flooding under multiple  
compound drivers. *Scientific Reports*, 12(1). <https://doi.org/10.1038/s41598-022-09652-5>, 2022.
- Ferrarin, C., Valentini, A., Vodopivec, M., Klaric, D., Massaro, G., Bajo, M., De Pascalis, F., Fadini, A., Ghezzi, M.,  
535 Menegon, S., Bressan, L., Unguendoli, S., Fettich, A., Jerman, J., Licer, M., Fustar, L., Papa, A., and Carraro, E.: Integrated  
sea storm management strategy: the 29 October 2018 event in the Adriatic Sea. *Nat. Hazards Earth Syst. Sci.*, 20, 73–93.  
<https://doi.org/10.5194/nhess-20-73-2020>, 2020.

- Giesen, R., Clementi, E., Bajo, M., Federico, I., Stoffelen, A., and Santoleri, R.: Copernicus Marine Service Ocean State  
540 Report, Issue 5. Section 4.3: The November 2019 record high water level in Venice, Italy. *Journal of Operational  
Oceanography*, 1–185. <https://doi.org/10.1080/1755876X.2021.1946240>, 2021.
- Gilleland, E., Ahijevych, D., Brown, B. G., Casati, B., and Ebert, E. E.: Intercomparison of spatial forecast verification  
545 methods. *Weather and Forecasting*, 24(5), 1416–1430. <https://doi.org/10.1175/2009WAF2222269.1>, 2009.
- Gumuscu, I., Islek, F., Yuksel, Y., and Sahin, C.: Spatiotemporal long-term wind and storm characteristics over the eastern  
Mediterranean Sea. *Regional Studies in Marine Science*, 63. <https://doi.org/10.1016/j.rsma.2023.102996>, 2023.
- Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., Nicolas, J., Peubey, C., Radu, R., Schepers,  
550 D., Simmons, A., Soci, C., Abdalla, S., Abellan, X., Balsamo, G., Bechtold, P., Biavati, G., Bidlot, J., Bonavita, M., ...  
Thépaut, J. N.: The ERA5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, 146(730), 1999–2049.  
<https://doi.org/10.1002/qj.3803>, 2020.
- Hervouet, J.-M., and Bates, P.: The TELEMAC modelling system Special issue. *Hydrological Processes*, 14(13), 2207–2208.  
555 [https://doi.org/10.1002/1099-1085\(200009\)14:13<2207::aid-hyp22>3.0.co;2-b](https://doi.org/10.1002/1099-1085(200009)14:13<2207::aid-hyp22>3.0.co;2-b), 2000.
- Hetzl, Y., Janekovic, I., and Pattiaratchi, C.: Assessing the ability of storm surge models to simulate coastal trapped waves  
around Australia. *Coasts and Ports*. <https://doi.org/10.3316/informit.929951406285439>, 2017.
- Lionello, P., Barriopedro, D., Ferrarin, C., Nicholls, R. J., Orlic, M., Raicich, F., Reale, M., Umgiesser, G., Vousdoukas, M.,  
560 and Zanchettin, D.: Extreme floods in Venice: characteristics, dynamics, past and future evolution (review article). *Nat.  
Hazards Earth Syst. Sci.*, 21, 2705–2731. <https://doi.org/10.5194/nhess-21-2705-2021>, 2021.
- Lionello, P., Cavaleri, L., Nissen, K. M., Pino, C., Raicich, F., and Ulbrich, U. (2012). Severe marine storms in the Northern  
565 Adriatic: Characteristics and trends. *Physics and Chemistry of the Earth*, 93–105. <https://doi.org/10.1016/j.pce.2010.10.002>
- Lionello, P., Galati, M. B., and Elvini, E.: Extreme storm surge and wind wave climate scenario simulations at the Venetian  
littoral. *Physics and Chemistry of the Earth*, 40–41, 86–92. <https://doi.org/10.1016/j.pce.2010.04.001>, 2010.
- Lovato, T., Androsov, A., Romanenkov, D., and Rubino, A.: The tidal and wind induced hydrodynamics of the composite  
570 system Adriatic Sea/Lagoon of Venice. *Continental Shelf Research*, 30(6), 692–706. <https://doi.org/10.1016/j.csr.2010.01.005>,  
2010.



- Luettich, R. A., Westerink, J. J., and Scheffner, N. W.: ADCIRC: an advanced three-dimensional circulation model for shelves, coasts, and estuaries. Report 1, Theory and methodology of ADCIRC-2DD1 and ADCIRC-3DL, 1992.
- 575
- Lyard, F., Allain, D. J., Cancet, M., Carrère, L., and Picot, N.: FES2014 global ocean tide atlas: design and performance. *Ocean Sci.*, 17, 615–649. <https://doi.org/10.5194/os-17-615-2021>, 2021.
- Maicu, F., Alessandri, J., Pinardi, N., Verri, G., Umgiesser, G., Lovo, S., Turolla, S., Paccagnella, T., and Valentini, A.:  
580 Downscaling with an unstructured coastal-ocean model to the Goro Lagoon and the Po River Delta branches. *Frontiers in Marine Science*, 8. <https://doi.org/10.3389/fmars.2021.647781>, 2021.
- Medugorac, I., Orlić, M., Janeković, I., Pasaric, Z., and Pasaric, M.: Adriatic storm surges and related cross-basin sea-level slope. *Journal of Marine Systems*. <https://doi.org/10.1016/j.jmarsys.2018.02.005>, 2018.
- 585
- Mentaschi, L., Besio, G., Cassola, F., and Mazzino, A.: Problems in RMSE-based wave model validations. *Ocean Modelling*, 72, 53–58. <https://doi.org/10.1016/j.ocemod.2013.08.003>, 2013.
- Mentaschi, L., Besio, G., Cassola, F., and Mazzino, A.: Performance evaluation of Wavewatch III in the Mediterranean Sea.  
590 *Ocean Modelling*, 90, 82–94. <https://doi.org/10.1016/j.ocemod.2015.04.003>, 2015.
- Mentaschi, L., Vousedoukas, M. I., García-Sánchez, G., Fernández-Montblanc, T., Roland, A., Voukouvalas, E., Federico, I., Abdolali, A., Zhang, Y. J., and Feyen, L.: A global unstructured, coupled, high-resolution hindcast of waves and storm surge. *Frontiers in Marine Science*, 10. <https://doi.org/10.3389/fmars.2023.1233679>, 2023.
- 595
- Mentaschi, L., Vousedoukas, M., Montblanc, T. F., Kakoulaki, G., Voukouvalas, E., Besio, G., and Salamon, P.: Assessment of global wave models on regular and unstructured grids using the Unresolved Obstacles Source Term. *Ocean Dynamics*, 70(11), 1475–1483. <https://doi.org/10.1007/s10236-020-01410-3>, 2020.
- 600
- Merrifield, M. A., Genz, A. S., Kontoes, C. P., and Marra, J. J.: Annual maximum water levels from tide gauges: Contributing factors and geographic patterns. *Journal of Geophysical Research: Oceans*, 118(5), 2535–2546. <https://doi.org/10.1002/jgrc.20173>, 2013.

- 605 Micaletto, G., Barletta, I., Mocavero, S., Federico, I., Epicoco, I., Verri, G., Coppini, G., Schiano, P., Aloisio, G., and Pinardi,  
N.: Parallel Implementation of the SHYFEM Model. *Geosci. Model Dev.*, 15, 6025–6046. <https://doi.org/10.5194/gmd-2021-319>, 2022.
- Mittermaier, M. P.: A strategy for verifying near-convection-resolving model forecasts at observing sites. *Weather and Forecasting*, 29(2), 185–204. <https://doi.org/10.1175/WAF-D-12-00075.1>, 2014.
- 610 Muis, S., Verlaan, M., Winsemius, H. C., Aerts, J. C. J. H., and Ward, P. J.: A global reanalysis of storm surges and extreme sea levels. *Nature Communications*, 7. <https://doi.org/10.1038/ncomms11969>, 2016.
- Muñoz, D. F., Yin, D., Bakhtyar, R., Mofstakhari, H., Xue, Z., Mandli, K., and Ferreira, C.: Inter-Model Comparison of Delft3D-FM and 2D HEC-RAS for Total Water Level Prediction in Coastal to Inland Transition Zones. *Journal of the American Water Resources Association*, 58(1), 34–49. <https://doi.org/10.1111/1752-1688.12952>, 2022.
- 615 Park, K., Federico, I., Di Lorenzo, E., Ezer, T., Cobb, K. M., Pinardi, N., and Coppini, G.: The contribution of hurricane remote ocean forcing to storm surge along the Southeastern U.S. coast. *Coastal Engineering*, 173. <https://doi.org/10.1016/j.coastaleng.2022.104098>, 2022.
- 620 Orlić, M., Kuzmić, M., and Pasarić, Z.: Response of the Adriatic Sea to the bora and sirocco forcing. *Continental Shelf Research*, 14, 91-116, 1994.
- 625 Pawlowicz, R., Beardsley, B., and Lentz, S.: Classical harmonic analysis including error estimates in MATLAB using T\_TIDE. *Computers and Geosciences*, 28, 929–937, 2022.
- Petti, M., Pascolo, S., Bosa, S., Bezzi, A., and Fontolan, G.: Tidal flats morphodynamics: A new conceptual model to predict their evolution over a medium-long period. *Water (Switzerland)*, 11(6). <https://doi.org/10.3390/w11061176>, 2019.
- 630 Pineau-Guillou, L., Ardhuin, F., Bouin, M. N., Redelsperger, J. L., Chapron, B., Bidlot, J. R., and Quilfen, Y.: Strong winds in a coupled wave–atmosphere model during a North Atlantic storm event: evaluation against observations. *Quarterly Journal of the Royal Meteorological Society*, 144(711), 317–332. <https://doi.org/10.1002/qj.3205>, 2018.
- 635 Pirazzoli, P. A., and Tomasin, A.: Recent evolution of surge-related events in the Northern Adriatic Sea. *Journal of Coastal Research*, 18, 537–554, 2022.

- Pringle, W. J., Wirasaet, D., Roberts, K. J., and Westerink, J. J.: Global storm tide modeling with ADCIRC v55: unstructured mesh design and performance. *Geosci. Model. Dev.*, 14, 1125–1145. <https://doi.org/10.5194/gmd-14-1125-2021>, 2021.
- 640 Raicich, F.: The sea level time series of Trieste, Molo Sartorio, Italy (1869-2021). *Earth System Science Data*, 15, 1749–1763. <https://doi.org/10.5194/essd-15-1749-2023>, 2023.
- Reimann, L., Vafeidis, A. T., Brown, S., Hinkel, J., and Tol, R.: Mediterranean UNESCO World Heritage at risk from coastal flooding and erosion due to sea-level rise. *Nature Communication*, 9:4161. <https://doi.org/10.1038/s41467-018-06645-9>, 2018.
- 645 Roberts, K. J., Pringle, W. J., and Westerink, J. J.: OceanMesh2D 1.0: MATLAB- based software for two-dimensional unstructured mesh generation in coastal ocean modeling. *Geosci. Model. Dev.*, 12, 1847–1868. <https://doi.org/10.5194/gmd-12-1847-2019>, 2019.
- 650 Roberts, N. M., and Lean, H. W.: Scale-selective verification of rainfall accumulations from high-resolution forecasts of convective events. *Monthly Weather Review*, 136(1), 78–97. <https://doi.org/10.1175/2007MWR2123.1>, 2008.
- Saha, S., Moorthi, S., Pan, H. L., Wu, X., Wang, J., Nadiga, S., Tripp, P., Kistler, R., Woollen, J., Behringer, D., Liu, H., Stokes, D., Grumbine, R., Gayno, G., Wang, J., Hou, Y. T., Chuang, H. Y., Juang, H. M. H., Sela, J., ... Goldberg, M.: The NCEP climate forecast system reanalysis. *Bulletin of the American Meteorological Society*, 91(8), 1015–1057. <https://doi.org/10.1175/2010BAMS3001.1>, 2010.
- 655 Saillour, T., Cozzuto, G., Ligorio, F., Lupoi, G., and Bourban, S. E.: Modeling the world oceans with TELEMAC. 2020 TELEMAC-MASCARET User Conference, 86–91, 2021.
- 660 Skok, G., and Roberts, N.: Analysis of Fractions Skill Score properties for random precipitation fields and ECMWF forecasts. *Quarterly Journal of the Royal Meteorological Society*, 142(700), 2599–2610. <https://doi.org/10.1002/qj.2849>, 2016.
- 665 Staneva, J., Wahle, K., Koch, W., Behrens, A., Fenoglio-Marc, L., and Stanev, E.: Coastal flooding: impact of waves on storm surge during extremes - a case study for the German Bight. *Nat. Hazards Earth Syst. Sci.*, 16, 2373–2389. <https://doi.org/10.5194/nhess-16-2373-2016>, 2016.
- Tiggeloven, T., Couasnon, A., van Straaten, C., Muis, S., and Ward, P. J.: Exploring deep learning capabilities for surge predictions in coastal areas. *Scientific Reports*, 11(1). <https://doi.org/10.1038/s41598-021-96674-0>, 2021.
- 670

- Toomey, T., Amores, A., Marcos, M., and Orfila, A.: Coastal sea levels and wind-waves in the Mediterranean Sea since 1950 from a high-resolution ocean reanalysis. *Frontiers in Marine Science*, 9. <https://doi.org/10.3389/fmars.2022.991504>, 2022.
- 675 Trotta, F., Fenu, E., Pinardi, N., Bruciaferri, D., Giacomelli, L., Federico, I., and Coppini, G.: A structured and unstructured grid relocatable ocean platform for forecasting (SURF). *Deep-Sea Research II*, 133, 54–75. <https://doi.org/10.1016/j.dsr2.2016.05.004>, 2016.
- Umgiesser, G., and Coauthors: The prediction of floods in Venice: Methods, models and uncertainty (review article). *Nat. Hazards Earth Syst. Sci.*, 21, 2679–2704. <https://doi.org/10.5194/nhess-21-2679-2021>, 2021.
- 680 Umgiesser, G., Canu, D. M., Cucco, A., and Solidoro, C.: A finite element model for the Venice lagoon: Development, set up, calibration and validation. *J. Marine Syst.*, 123–145. <https://doi.org/10.1016/j.jmarsys.2004.05.009>, 2004.
- 685 Vannucchi, V., Taddei, S., Capecchi, V., Bendoni, M., and Brandini, C.: Dynamical downscaling of era5 data on the north-western mediterranean sea: From atmosphere to high-resolution coastal wave climate. *Journal of Marine Science and Engineering*, 9(2), 1–29. <https://doi.org/10.3390/jmse9020208>, 2021.
- Vousdoukas, M. I., Clarke, J., Ranasinghe, R., Reimann, L., Khalaf, N., Duong, T. M., Ouweneel, B., Sabour, S., Iles, C. E., 690 Trisos, C. H., Feyen, L., Mentaschi, L., and Simpson, N. P.: African heritage sites threatened as sea-level rise accelerates. *Nature Climate Change*, 12(3), 256–262. <https://doi.org/10.1038/s41558-022-01280-1>, 2022.
- Vousdoukas, M. I., Mentaschi, L., Voukouvalas, E., Verlaan, M., Jevrejeva, S., Jackson, L. P., and Feyen, L.: Global probabilistic projections of extreme sea levels show intensification of coastal flood hazard. *Nature Communications*, 9(1). 695 <https://doi.org/10.1038/s41467-018-04692-w>, 2018.
- Wang, X., Verlaan, M., Veenstra, J., and Lin, H. X.: Data-assimilation-based parameter estimation of bathymetry and bottom friction coefficient to improve coastal accuracy in a global tide model. *Ocean Sci.*, 18, 881–904. <https://doi.org/10.5194/os-18-881-2022>, 2022.
- 700 Weatherall, P., Marks, K. M., Jakobsson, M., Schmitt, T., Tani, S., Arndt, J. E., Rovere, M., Chayes, D., Ferrini, V., and Wigley, R.: A new digital bathymetric model of the world's oceans. *Earth and Space Science*, 2(8), 331–345. <https://doi.org/10.1002/2015EA000107>, 2015.

- 705 Weisberg, R. H., and Zheng, L.: Hurricane storm surge simulations comparing three-dimensional with two-dimensional formulations based on an Ivan-like storm over the Tampa Bay, Florida region. *Journal of Geophysical Research: Oceans*, 113(12). <https://doi.org/10.1029/2008JC005115>, 2008.
- World Meteorological Organization: Guide to storm surge forecasting, 2011.
- 710 Ye, F., Zhang, Y., Yu, H., Sun, W., Moghimi, S., Myers, E., Nunez, K., Zhang, R., Wang, H., Roland, A., Martins, K., Bertin, X., Du, J., and Liu, Z.: Simulating storm surge and compound flooding events with a creek-to-ocean model: Importance of baroclinic effects. *Ocean Modelling*, 145. <https://doi.org/10.1016/j.ocemod.2019.101526>, 2020.
- 715 Yu, C. S., Decouttere, C., and Berlamont, J.: Storm surge simulations in the Adriatic Sea. In G. Gambolati (Ed.), *CENAS. Coastline Evolution of the Upper Adriatic Sea due to Sea Level Rise and Natural and Anthropogenic Land Subsidence* (pp. 207–232), 1998.
- Zaggia, L., Lorenzetti, G., Manfè, G., Scarpa, G. M., Molinaroli, E., Parnell, K. E., Rapaglia, J. P., Gionta, M., and Soomere, T.: Fast shoreline erosion induced by ship wakes in a coastal lagoon: Field evidence and remote sensing analysis. *PLoS ONE*, 12(10). <https://doi.org/10.1371/journal.pone.0187210>, 2017.
- Zampato, L., Umgiesser, G., and Zecchetto, S.: Storm surge in the Adriatic Sea: observational and numerical diagnosis of an extreme event. *Advances in Geosciences*, 7, 371–378, 2006.
- 725 Zampato, L., Umgiesser, G., and Zecchetto, S.: Sea level forecasting in Venice through high resolution meteorological fields. *Estuarine, Coastal and Shelf Science*, 75, 223–235. <https://doi.org/10.1016/j.ecss.2007.02.024>, 2007.
- Zhang, Y., and Baptista, A. M.: SELFE: A semi-implicit Eulerian–Lagrangian finite-element model for cross-scale ocean circulation. *Ocean Model.*, 21, 71–96. <https://doi.org/10.1016/j.ocemod.2007.11.005>, 2008.
- 730 Zhang, Y. J., Fernandez-Montblanc, T., Pringle, W. J., Yu, H. C., and Cui, L.: Global seamless tidal simulation using a 3D unstructured-grid model (SCHISM v5.10.0). *Geosci. Model. Dev.* <https://doi.org/10.5194/gmd0-16-2565-2023>, 2023.
- 735 Zhang, Y., Ye, F., Stanev, E. V., and Grashorn, S.: Seamless cross-scale modeling with SCHISM. *Ocean Model.*, 102, 64–81. <https://doi.org/10.1016/j.ocemod.2016.05.002>, 2016.

Zingerle, C., and Nurmi, P.: Monitoring and verifying cloud forecasts originating from operational numerical models. *Meteorological Applications*, 15(3), 325–330. <https://doi.org/10.1002/met.73>, 2008.

740