

All review comments and responses to

“Towards using state-of-the-art climate models to help constrain estimates of unprecedented UK storm surges”

By Tom Howard and Simon David Paul Williams

History of the review (taken from the NHESS website)

RC1: 'Comment on nhess-2021-184', Andreas Sterl, 13 Aug 2021 

└ RC2: 'Reply on RC1', Andreas Sterl, 13 Aug 2021 


└ AC1: 'Reply on RC2', Tom Howard, 18 Aug 2021 


└ AC2: 'Reply on AC1', Tom Howard, 18 Aug 2021 

RC3: 'Comment on nhess-2021-184', Jonathan Tawn, 27 Aug 2021 

└ AC3: 'Reply on RC3', Tom Howard, 09 Sep 2021 

└ AC4: 'Reply on AC3', Tom Howard, 09 Sep 2021 

└ AC5: 'Reply on AC4', Tom Howard, 09 Sep 2021 

└ RC4: 'Reply on AC5', Jonathan Tawn, 16 Sep 2021 

└ AC6: 'Reply on RC4', Tom Howard, 18 Sep 2021 

RC1:

Review of

Towards using state-of-the-art climate models to help constrain estimates of unprecedented UK storm surges by T. Howard and S. Williams

Recommendation

Minor / major revisions - depend on editor's decision on investigating a Gumbel fit (see Discussion section below).

Synopsis

The paper investigates whether it is possible to exploit long climate model runs to estimate long return times of surges along the UK coast. To this end, the output of a 483-year integration of HadGEM3-GC3-MM is used to force the CS3 barotropic surge model. CS3 results are then compared to the (much shorter) observational records. Extreme Value Theory is used to infer long return times. A great deal of effort is put into estimation of the parameters of a GEV fitted to the model results. The shape parameter of a GEV distribution is found to be the source of the largest uncertainty in the return-level estimate.

The paper mainly deals with the skew surge because this measure has been shown to be independent of the tidal phase. In the case of Shearness, the modelling results contradict this view. The skew surge created by a specific wind storm differs by 20%, depending on whether the wind storm occurred during neap tide or during spring tide, respectively.

Discussion

Coastal protection works (dykes, flood barriers, etc) should protect against events of a magnitude that has not been reached during the observational record. Typically, such works are designed to protect against a water level occurring once in 1000 or 10,000 years. This time is much longer than the length of the observational records, which are usually not longer than 100 years. Extrapolating over orders of magnitude naturally leads to large uncertainties. Replacing observations by model-generated series, which often are much longer than the observational record, is therefore a good idea. The approach has been introduced nearly 20 years ago by Van Den Brink et al. (2004): Improving 10⁴-year surge level estimates using data of the ECMWF seasonal prediction system. *Geophys. Res. Lett.*, 31, L17210, doi: 10.1029/2004GL020610. This fact should be acknowledged in the present paper.

The authors fit the three parameter (location, shape, scale) GEV distribution to the model results and find that the largest uncertainty comes from the scale parameter. That the scale parameter is the most uncertain parameter in a GEV fit is well known. It is most heavily determined by the highest values in the annual-maximum series. A way around is to use the two parameter (location and scale) Gumbel distribution. Especially for long time series, it often gives superior results. Looking at Fig. 2c and d, a value of zero for the scale parameters is not inconsistent with at least the observational estimates. Whether the Gumbel distribution is a good approximation can be tested by the procedure explained in Van den Brink and Können (2011): Estimating 10000-year return values from short time series. *Int. J. Climatol.*, 31:115-126, doi: 10.1002/joc.2047.

I am aware of the fact that it would require a lot of work to calculate the Gumbel-fits and test for its suitability. However, I feel that the results shown in the paper could be much improved by eliminating the uncertainty that is inherent in the scale parameter. The editor should decide whether (s)he requires this analysis to be done.

Detailed comments

The paper is very well written, and I have only a few minor comments.

p 12, l 300/301 Discussing Fig. 2c you refer to the “spread of the shape parameters diagnosed [...] from the simulation”, but I cannot see any spread in the simulation-derived shape parameters. Spread is only depicted for the CFB estimates. Please clarify.

sec. 6.1 / Fig. 6 You introduce a kernel. How did you obtain it? A short explanation of the procedure would be helpful. I am also not sure about the purpose of the kernel. Is it to extend the length of an episode, or its magnitude?

RC2:

I obviously confused 'shape' and 'scale' in my comment. In the second para of the Discussion section it is the scale parameter that causes the largest uncertainty and that seems to be compatible with zero. Sorry for any confusion that may have arisen..

AC1:

Response to review by Andreas Sterl of *Towards using state-of-the-art climate models to help constrain estimates of unprecedented UK storm surges.*

by T. Howard and S. Williams, submitted 2021

(henceforth HW21)

This response: Tom Howard Aug 2021

Thank you for your review and for raising a very interesting question. I have amended the draft to address the minor points, and included a discussion of your main point further below. I will attempt to upload the amended draft to accompany this response.

Regarding the minor points first (quotes from the review are shown in red).

The approach has been introduced nearly 20 years ago by Van Den Brink et al. (2004): Improving 10^4 year surge level estimates using data of the ECMWF seasonal prediction system. *Geophys. Res. Lett.*, 31, L17210, doi: 10.1029/2004GL020610. This fact should be acknowledged in the present paper.

I apologise for failing to cite this very relevant paper in the first draft. I have now acknowledged this contribution.

Discussing Fig. 2c you refer to the “spread of the shape parameters diagnosed [...] from the simulation”, but I cannot see any spread in the simulation-derived shape parameters. Spread is only depicted for the CFB estimates. Please clarify.

Thank you for pointing this out. I was referring to the spread associated with the spatial variations, rather than the uncertainty. I have clarified this in the paper as follows:

“The spread (i.e. the size of the spatial variations) of the shape parameters...”

You introduce a kernel. How did you obtain it? A short explanation of the procedure would be helpful. I am also not sure about the purpose of the kernel. Is it to

extend the length of an episode, or its magnitude?

I have added the following sentences explaining the choice and purpose of the kernel:

“The kernel was designed to represent the important features of the RACMO-driven surge-only simulation, i.e., the approximate duration and shape of the time series plot. The purpose of convolution with the kernel is to identify those events which correlate well (in terms of their time series plot) with the RACMO-driven simulation, in other words, events which not only produce a large surge, but are also of comparable duration to the RACMO-driven simulation. The kernel was not used to modify events, but simply to identify significant ones.”

Turning now to your main point regarding whether to fix the shape parameter at zero. This is prompted by the two papers by van den Brink and Können:

van den Brink and Können (2011): Estimating 10000-year return values from short time series. *Int. J. Climatol.*, 31:115-126, doi: 10.1002/joc.2047 , and their related 2008 paper (both referred to here as vdB&K).

Thank you for reminding me of vdB&K’s very interesting approach. I spent a long time (back in 2017) studying their 2008 paper, but it is only now, in testing the method on our own data in order to complete this response, that I am beginning to understand it.

HW21 was rooted in the methodology of CFB2018, which forms the current UK guidance on sea level extremes:

https://assets.publishing.service.gov.uk/media/603652cce90e0740b7caac9d/Coastal_flood_boundary_conditions_for_the_UK_2018_update_-_technical_report.pdf

(full reference in HW21). CFB2018 was developed in consultation with Professor Jonathan Tawn:

<https://www.maths.lancs.ac.uk/~tawn/>

Jonathan advised me not to fix the shape parameter as this gives false confidence intervals. The position is laid out in the textbook by Stuart Coles (full reference in HW21). Coles provides an illustration of a case similar to the case in section 4.5 of vdB&K (2011), and Coles discusses it as follows (Coles page 64)<Quote>:

Reduction of uncertainty is desirable, so that if the Gumbel model could be trusted its inferences would be preferred. But can the model be trusted? The extremal types theorem provides support for modelling block maxima with the GEV family of which the Gumbel family is a subset. The data [in Coles's example 3.4.1, which, like vdB&K(2011) example 4.5, has an estimated shape parameter close to zero] suggest that a Gumbel model is plausible, but this does not imply that other models are not. Indeed, the maximum likelihood estimate within the GEV family is not in the Gumbel family (although, in the sense that the estimated shape parameter is close to zero, it is "close"). There is no common agreement about this issue, but the safest option is to accept there is uncertainty about the value of the shape parameter --- and hence whether the Gumbel model is correct or not --- and to prefer the inference based on the GEV model. The larger measures of uncertainty generated by the GEV model then provide a more realistic quantification of genuine uncertainties involved in model extrapolation.

<End Quote> (The red highlighting is mine)

On the other hand, we know that unconstrained GEV fits to short record lengths can give implausible shape parameters. One way to fix this is to put a prior on the shape parameter (see for example Martins and Stedinger 2000, full reference in HW21). This was the approach used in CFB2018 (again under advice from Jonathan Tawn).

In HW21 we find that the spatial variations in the shape parameter diagnosed from the tide gauges are also seen in the shape parameter as diagnosed from the simulation. We have argued that this finding supports the credibility of the spatial variations. Thus, it would be inconsistent to adopt an approach of fixing the shape parameter at zero within HW21, but for completeness I have tested a Gumbel fit to the annual maxima following the method of vdB&K and I show the results here.

I have used our data to make plots analogous to vdB&K 2008 Figure 3, i.e., their $\Delta\hat{X}_n$ vs. Gumbel variate plot, comparing Gumbel fits to the annual maxima with the fits used in our paper and in CFB2018. Fig. 1 is the plot for the observations:

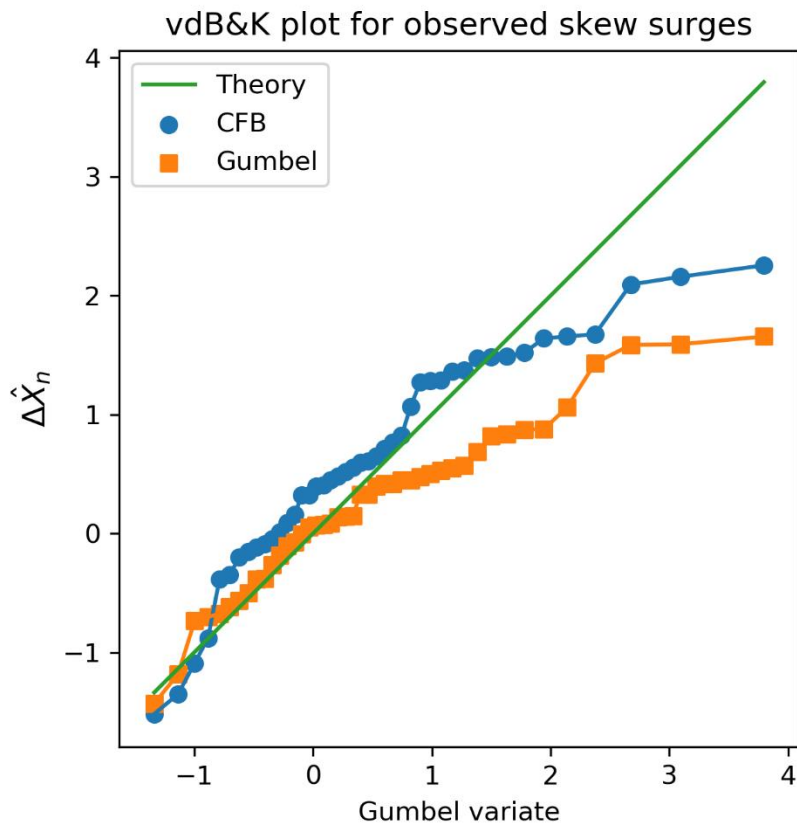


Fig. 1. Each data point represents a UK tide gauge. The points labelled “CFB” show $\Delta\hat{X}_n$ where \hat{F} is determined by the CFB2018 method (constrained GPD fit to peaks over a threshold, as described in HW21), and the points labelled “Gumbel” show $\Delta\hat{X}_n$ where \hat{F} is determined by a Gumbel fit to the annual maxima.

I can't see strong evidence here that the Gumbel fit is better overall. I was a bit concerned that I had not been rigorous about checking for independence of the events shown. This would not be straightforward to do with my current software setup, but a simple first fix is to miss out closely-neighbouring tide gauges. The following plots show results from every second tide gauge (Separation=2), every third tide gauge (Separation=3), etc. (Separation=1 is just a duplicate of the above plot).

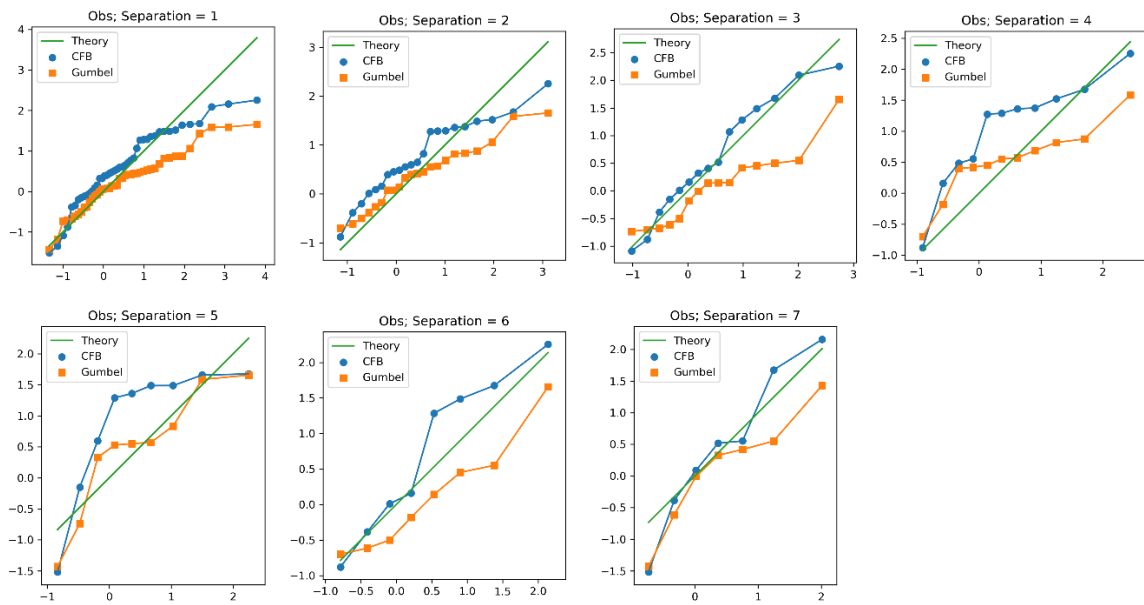


Fig. 2. See main text.

Again, there does not appear to be much support for preferring the Gumbel fit.

(continued...)

I followed the same procedure for the simulated skew surges to make Figs 3 and 4.

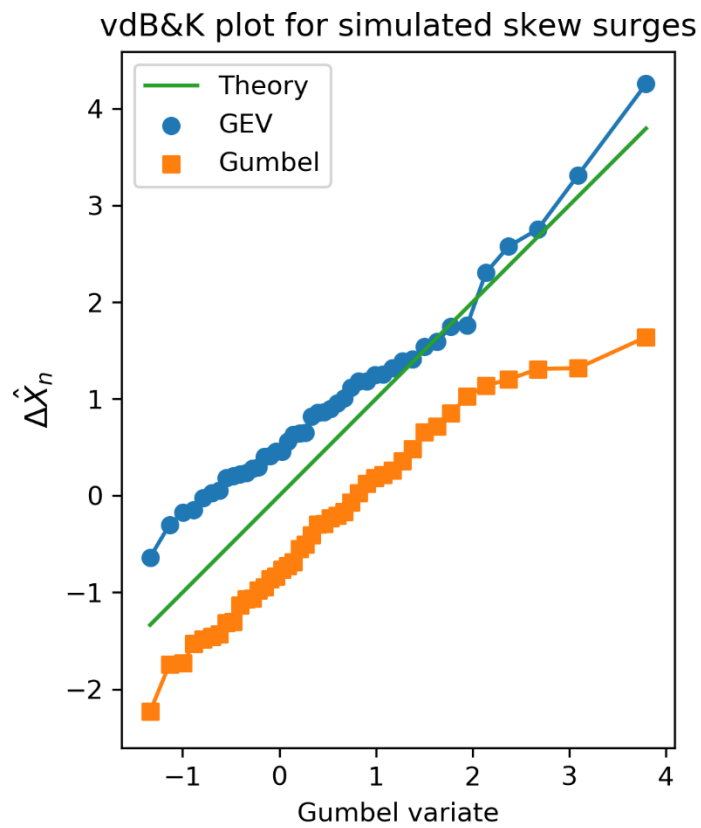


Fig. 3. As Fig. 1, but for the simulated skew surges. Each data point represents a UK tide gauge. The points labelled “GEV” show $\Delta \hat{X}_n$ where \hat{F} is determined by a GEV fit to the annual maxima, and the points labelled “Gumbel” show $\Delta \hat{X}_n$ where \hat{F} is determined by a Gumbel fit to the annual maxima.

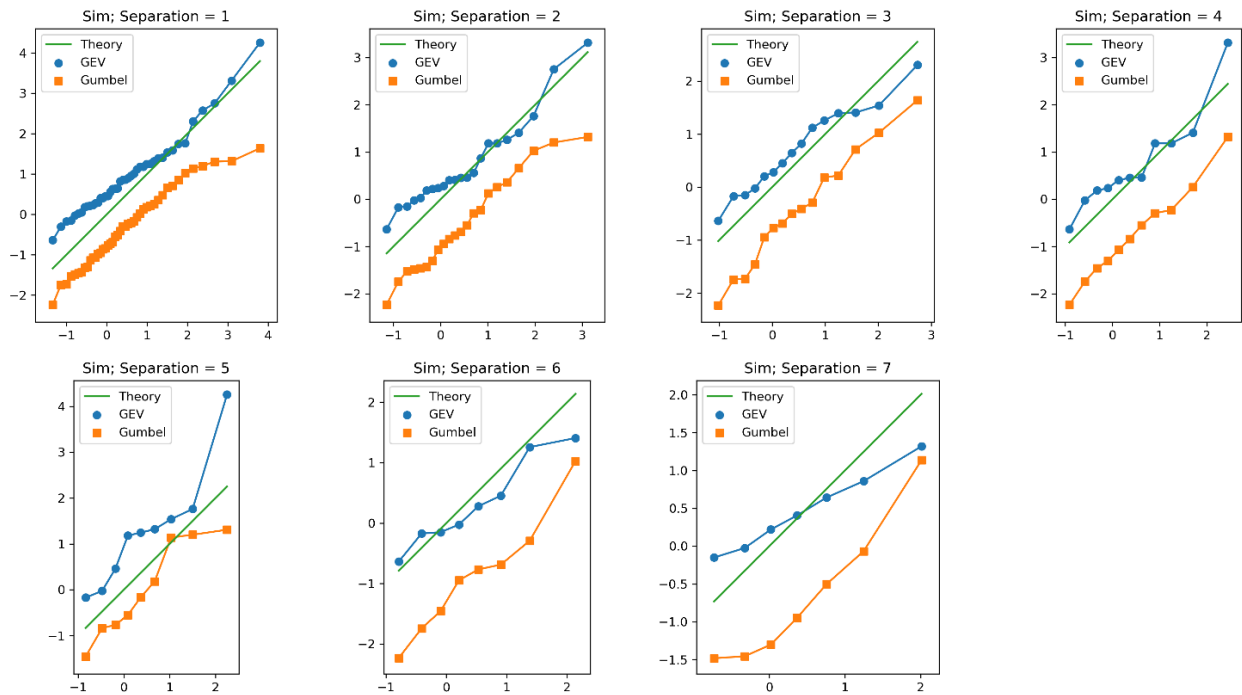


Fig. 4. As Fig. 2, but for the simulated skew surges.

Again, this does not seem to me to give support for preferring the Gumbel fit. Even if it did, I would not think it defensible to fix the shape parameter at a single value (e.g., zero). Looked at from the point of view of applying a prior to the shape parameter, fixing it at a single value seems to be equivalent to asserting that we are sure that no other value is plausible, and we are sure that the shape parameter does not vary by location. I cannot support either of those assertions.

Note to self: internal reference for figures: fig_KK

The reviewer responded with a personal communication which I have copied to the editor.

AC2: This was a diff.pdf between the original submission and the first revision. There seems little point in including it here because we later have the diff between the original submission and the second revision.

RC3:

Discussion of Towards using state-of-the-art climate models to help constrain estimates of unprecedented UK storm surges

by Tom Howard and Simon David Paul Williams

Discussants: Eleanor D'Arcy and Jonathan Tawn

August 27, 2021

Overview

This paper proposes using climate model simulations to aid with still water level return level estimation and address problems that arise with the statistical approach when tide gauges have short record lengths. The model they present is the HadGEM3-GC3-MM to generate a dataset of 483-year present-day storm surges at sites on the UK tide gauge network. They compare the skew surge simulations to using only observations when fitting extreme value models by estimating parameters. The spatial distribution of the parameters for each dataset are generally well correlated. However, there is a negative bias in the simulation approach in the shape parameter estimate of the generalised Pareto distribution (GPD); the authors discuss this in detail and suspect it is due to a pitfall in the shelf sea model (CS3).

The paper also investigates the interaction between skew surge and peak tide at Sheerness. They study the effect that changes in timing between atmospheric forcings and tide has on skew surge. Additionally, they review the independence assumption of skew surge and peak tide used in the JPM. Using their model simulations, they show that extreme skew surges are more likely to occur on neap tide - this agrees with the results of Williams et al. (2016) (supplementary material).

General Comments

The paper is well written, in the proceeding sections we have discussed parts of the paper that pose interesting areas for future research and noted some technical corrections. The results are well presented and the figures well explained, giving a clear justification for the proposed model. The appendix provides strong support for ideas mentioned in the main paper. The authors have recognised potential downfalls with different aspects of the model and presented some initial investigation into these (for example, the discussion of why the shape parameter is more negative for simulations on pg. 15/16).

Specific Comments

Comparison of Model and Observed Data

Figure B1 shows some major departures between the model and observed data across the distribution of skew surges but particularly in the tails. In no sites is the model giving as high quantiles as the observed data. However, these departures have a systematic feature which is consistent over spatial regions, e.g., south-west and north-west UK. This suggests that it should be possible to account for these departures through a smooth spatial function which maps the

differences in quantiles between the observations and model data. With this adjustment it is possible that the currently identified under-estimation may be corrected before making the tail based GEV/GPD comparisons you draw.

Shape Parameters with Similar Spatial Patterns

One of your exciting findings is that the spatial pattern of the shape parameter estimates is similar for the CFB2018 estimates from observed data and your estimates from the model, but with a systematic bias between them. This suggests using an alternative to the CFB2018 approach by explicitly exploiting this finding. Let $\xi_{obs}(x)$ and $\xi_{model}(x)$ be the shape parameters for the observed and model data respectively, for all gauged sites x . What you are saying is that you believe in the spatial variation of $\xi_{model}(x)$ but not its mean value. So you believe that $\xi_{obs}(x) = \xi_{model}(x) + \xi_{bias}$, where ξ_{bias} is a fixed constant that does not vary over x . Your estimates indicate that $\xi_{bias} > 0$. Fixing estimates of $\xi_{model}(x)$ for all x and fitting the function of $\xi_{obs}(x)$ over x now means only one parameter, ξ_{bias} , needs estimating. This could lead to substantial reductions in the uncertainty of $\xi_{obs}(x)$ estimates over x .

Penalised Likelihood

The shape parameter estimates in Figure 2 (d) based on the 483 years of model data show some site-to-site variations which are more pronounced than the broader smooth variations across coastlines. This suggests that they would also benefit from the penalty-based approach used in CFB2018 work.

At a few points you state that the prior/penalty is subjective and that this is a disadvantage. Yet you also point that the smooth pattern of the shape parameter estimates this gives for the observed data agrees well with the similar unpenalised estimates using model data. You say this is a really positive feature of the model data, we would also take this as supporting the value of the process of creating penalised estimates from the observed data. The penalty is giving something meaningful, so the effect of the claimed “subjectivity” is positive for the observed data analysis. The prior that was selected in the CFB2018 work was not subjective in the traditional sense of a subjective prior in Bayesian methods. It was actually a data-based prior which corresponds to an empirical Bayesian prior, using all the information that separately estimated shape parameters for UK skew surge provide. The effect of this was simply to move shape parameter estimates more towards the UK average, with the larger changes coming for sites with shorter record lengths.

Investigation of Interaction between Tide and Skew Surge

We are really pleased to see the use of numerical models to explore systematically the widely claimed property that skew surges are independent of their associated peak tidal values. The work in Section 5 where you explore the timing effects of skew surge events relative to tides is very illuminating. It adds greatly to the existing empirical evidence for this feature at Sheerness. It would be interesting to have your opinions about what interactions are expected elsewhere in the UK given that the CFB2018 estimates over UK sites all assume independence. Since this analysis is based on simulations from HadGEM3-GC3-MM, this presents a physically-based justification for dependence between skew surge and peak tide at Sheerness. It suggests future research is required to investigate this. It may be that interaction occurs to some extent everywhere, but at a practical level it is not important apart from some locations.

In current work we have been investigating this feature empirically at a limited number of sites. Here we report some of the methods and findings using data from the tide gauges at Sheerness and

Heysham (located close to Workington which you analyse). We study data from Heysham in 1964-2016, of which 17.5% is missing, and 1980-2016 at Sheerness, where 9.1% is missing.

We define extreme skew surges as exceedances of the 0.95 quantile at each site. Figure 1 shows scatter plots of extreme skew surges against their associated ranked peak tide; these are ranked so that 1 corresponds to the smallest observation and n is the largest, where n is the total number of observations. If the two components were independent, we expect extreme skew surges to be uniformly distributed over ranks. We test this using a Kolmogorov-Smirnov test for uniformity; the p value at Heysham is 0.0224 whilst at Sheerness it is 1.40×10^{-7} . Clearly the p value at Sheerness is much smaller, and provides statistical evidence that extreme skew surges are not independent of peak tide. However, at Heysham, there is not sufficient evidence to reject the claim of independence at the 0.01 significance level. This is clear in Figure 1, where more extreme skew surges occur on lower tides at Sheerness - this agrees with your findings and those of Williams et al. (2016).

We investigate this further by looking at skew surge and peak tide dependence on a month-by-month basis, using ideas from Williams et al. (2016). Figure 2 compares the distribution of peak tides associated

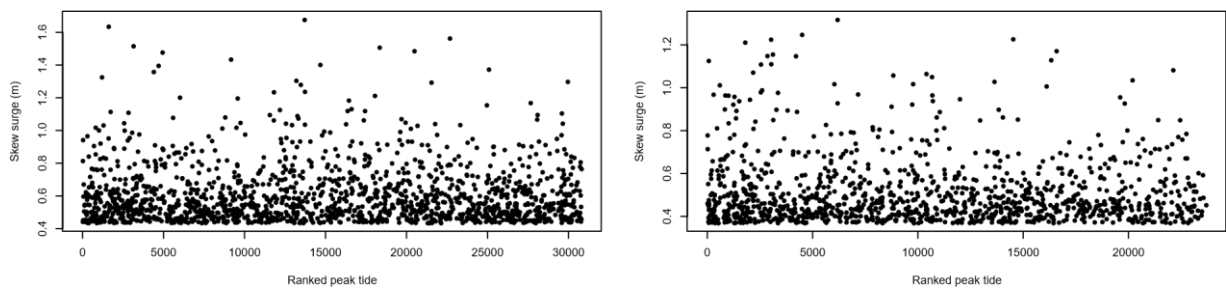


Figure 1: Extreme skew surge observations against ranked peak tide at Heysham (left) and Sheerness (right).

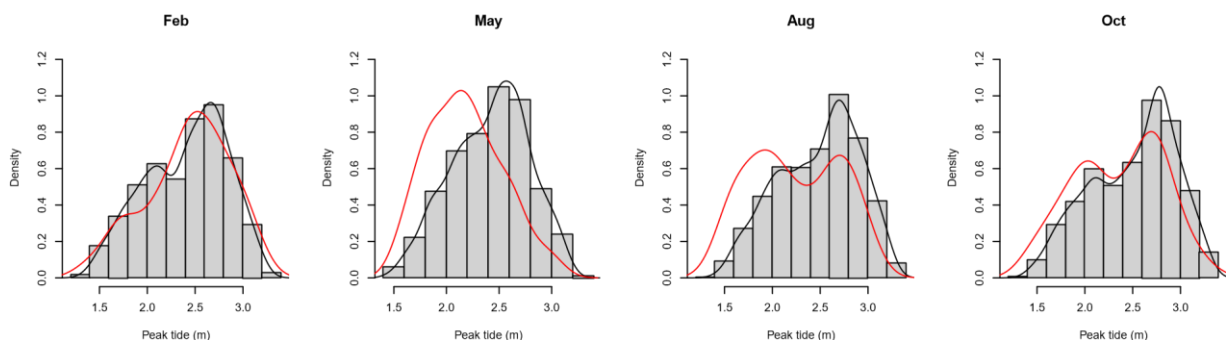


Figure 2: Monthly distributions of peak tides at Sheerness in February, May, August and October. The probability density function of all peak tides (black) and peak tides associated with extreme skew surge (red) are interpolated onto each distribution.

with all skew surges and the distribution of peak tides associated with extreme skew surges for February, May, August and October. If peak tide and skew surge are independent, these two distributions should be the same, up to sampling variation. We estimate the probability density function (pdf) using a Gaussian kernel density estimate, and use an Anderson-Darling test to check if peak tides come from the same distribution as peak tides associated with extreme skew surges. Figure 2 highlights how the dependence between skew surge and peak tide is changing with the time of year; the distribution of peak tides and the peak tides associated with extreme skew surge are

most different in May and least different in February. In May, the mode of the distribution of peak tides associated with extreme skew surges has shifted to a lower value than the distribution of all peak tides. Results from the Anderson Darling test for every month tell us there is insufficient evidence to reject the null hypothesis that peak tides and the peak tides associated with extreme skew surge come from the same distribution in February, March, September, November and December. In these months, we conclude it is reasonable to assume skew surge and peak tide are independent. In the remaining months, we find sufficient evidence to suggest the two components are dependent. We believe this poses a really interesting area for further research.

In D’Arcy et al. (2021) we fit a GPD (see Coles (2001) for details) to extreme skew surges that accounts for seasonal variations, since they are more extreme in the winter. Our results show this is an improvement on a standard GPD fit, where skew surges are assumed to be independent of peak tide. Here, we investigate adding a covariate of peak tide into our skew surge model to account for the dependence found at Sheerness. D’Arcy et al. (2021) defines extreme skew surges as exceedances of the monthly 0.95 quantile. Non-stationarity is accounted for in the scale parameter of the GPD through a daily covariate $d = 1, \dots, 365$:

$$\sigma_d = a + b \sin\left(\frac{2\pi}{365}(d - \phi)\right) \quad (3.1)$$

for $a, b, \phi \in \mathbb{R}$ parameters to be estimated. Note that the shape parameter of the GPD is fixed across months. To account for the dependence between skew surge and tide, we now also consider the following parameterisation on the scale parameter:

$$\sigma_{d,t} = a + b \sin\left(\frac{2\pi}{365}(d - \phi)\right) + ct \quad (3.2)$$

where $c \in \mathbb{R}$ is another parameter to be estimated, and t the associated peak tide observation. We fit a GPD with both (3.1) and (3.2) formulations to extreme skew surges at Heysham and Sheerness. The Akaike information criteria (AIC), frequently used for model selection, suggests that formulation (3.1) gives a better model fit at Heysham, i.e., an independence conclusion is supported by this analysis. Whereas, AIC suggests that the parametrisation in equation (3.2) yields a better fit at Sheerness. By ordering peak tide observations from smallest to largest and calculating $\sigma_{d,t}$ at its winter peak ($d = 365$) for each value, we observe a 15% reduction in the scale parameter as tides increase at Sheerness, which corresponds to an equal percentage reduction in the skew surge quantiles for excesses of the December skew surge threshold. We also compare the model fits at each site using a Likelihood Ratio Test. At Heysham the p value is 0.657 which provides insufficient evidence to reject the null hypothesis, that is the simpler model (in equation (3.1)) is sufficient. Whereas, at Sheerness we get a much smaller p value of 0.059, which provides statistically significant evidence at the 0.1 significance level to reject the null hypothesis and conclude that the more complex model is required. This highlights the importance of accounting for skew surge and peak tide dependence at sites where the independence assumption is not justified.

Ungauged Sites

We feel you do not do full justice to your developments given the focus is on comparisons made at gauged sites with long and trustworthy records. The real value in modelled data is the ability to give estimates at other sites. This could be the focus of a natural follow up paper.

Other Comments

1. In the abstract, you say “results suggest an event of this magnitude has an expected frequency of about 1 in 500 years at [Sheerness]” when referring to the North Sea floods of 1953. In the paper, the only result to show this is presented in Section 6.2: “the fact that the 483-year surge-only simulation produces more than one event of comparable magnitude to the simulated 1953 event suggest the return period of the 1953 atmospheric forcing is less than 483 years.” If this is the justification, it would be better to more formally quantify this using your statistical models.
2. In the Introduction (line 51) you list assumptions required to fit extreme value models as ‘events are effectively random and statistically independent of each other.’ But what about events being identically distributed (or stationary), it is clear that tide and storm surge (or skew surge) are not stationary as both process exhibit seasonality. Instead of “effectively random” it would be more mathematically correct to say “stochastic.” Extreme value methods also handle dependence, so “independent of each other” is not formally required.
3. On line 84 climate change is discussed. Tide gauge observations will exhibit an approximately linear mean trend due to sea level rise, it is unclear whether this has been removed before comparison with the simulations in Section 4. It is likely that removing this change will not change the results significantly, but it is important to remove this non-stationary effect before fitting extreme value models.
4. Figure C1 shows the reduction in uncertainty on shape parameter when record lengths are increased in the tide gauge network, but it would be nice to see this for more record lengths that are equally spaced (say 10 to 500 years in increments of 10 years) to really highlight this - with measures of uncertainty as in Figure C1 (d).

Technical Corrections

1. Table 1 gives a list of useful acronyms and symbols, it would be help to include what ‘GC3’ and ‘MM’ stand for in ‘HadGEM3-GC3-MM’ as it is not mentioned in text.
2. Equation (2) has a surplus close bracket.
3. Equation (6), should this derivative be evaluated at $L = \log(u)$ rather than $y = u$ since y is the return level (as in equation (1))?
4. Figures 2 (c) and (d) would benefit from having 95% confidence intervals on shape parameters using the model data.
5. Line 283: We think the wrong figure is referenced here.
6. Line 319: “They usedtime-series.” This has nothing to do with how CFB2018 estimate the shape parameters and should be cut. It only links to the derivation of SWL return levels.
7. Line 349: You hint that the reason for the disagreement between the methods could be due the short observed data series. But Figure B1 shows major departures between the observed data and the model data in the body of the distributions to a level that could in no way be due to short observed series and must be that the model data does not reproduce well the observed data. Linked to this, in discussing Figure 1 you say the agreement at both sites is “excellent”. With the amount of data at Sheerness it is clear that there are major disagreements that cannot be accounted by sampling variations.
8. Line 421: You suddenly mention a “kernel” to spread the duration of events but fail to provide any information. A little detail would be helpful here.

9. Figure E2: A description of the different points would be useful.

References

- Coles, S. G. (2001). *An Introduction to Statistical Modeling of Extreme Values*. Springer, London.
- D’Arcy, E., Tawn, J., Joly-Laugel, A., and Sifnioti, D. E. (2021). Accounting for seasonality in extreme sea level estimation (*in preparation*).
- Williams, J., Horsburgh, K. J., Williams, J. A., and Proctor, R. N. (2016). Tide and skew surge independence: New insights for flood risk. *Geophysical Research Letters*, 43(12):6410–6417.
-
-

AC3:

Tom Howard 9 Sep. 21

Author response to:

Discussion of Towards using state-of-the-art climate models to help constrain estimates of unprecedented UK storm surges

by Tom Howard and Simon David Paul Williams

Discussants: Eleanor D’Arcy and Jonathan Tawn

August 27, 2021

Thank you both very much for your thorough and interesting review. I have pasted in the whole discussion in dark orange font colour; my responses are shown in blue.

Overview

This paper proposes using climate model simulations to aid with still water level return level estimation and address problems that arise with the statistical approach when tide gauges have short record lengths. The model they present is the HadGEM3-GC3-MM to generate a dataset of 483-year present-day storm surges at sites on the UK tide gauge network. They compare the skew surge simulations to using only observations when fitting extreme value models by estimating parameters. The spatial distribution of the parameters for each dataset are generally well correlated. However, there is a negative bias in the simulation approach in the shape parameter estimate of the generalised Pareto distribution (GPD); the authors discuss this in detail and suspect it is due to a pitfall in the shelf sea model (CS3).

The paper also investigates the interaction between skew surge and peak tide at Sheerness. They study the effect that changes in timing between atmospheric forcings and tide has on skew surge. Additionally, they review the independence assumption of skew surge and peak tide used in the JPM. Using their model simulations, they

show that extreme skew surges are more likely to occur on neap tide - this agrees with the results of Williams et al. (2016) (supplementary material).

General Comments

The paper is well written, in the proceeding sections we have discussed parts of the paper that pose interesting areas for future research and noted some technical corrections. The results are well presented and the figures well explained, giving a clear justification for the proposed model. The appendix provides strong support for ideas mentioned in the main paper. The authors have recognised potential downfalls with different aspects of the model and presented some initial investigation into these (for example, the discussion of why the shape parameter is more negative for simulations on pg. 15/16).

Specific Comments

Comparison of Model and Observed Data

Figure B1 shows some major departures between the model and observed data across the distribution of skew surges but particularly in the tails. In no sites is the model giving as high quantiles as the observed data. However, these departures have a systematic feature which is consistent over spatial regions, e.g., south-west and north-west UK. This suggests that it should be possible to account for these departures through a smooth spatial function which maps the differences in quantiles between the observations and model data. With this adjustment it is possible that the currently identified under-estimation may be corrected before making the tail based GEV/GPD comparisons you draw.

Have added the following to the description accompanying Figure B1:

“Figure B1 shows some major departures between the model and observed data across the distribution of skew surges, but particularly in the tails. The model does not give higher quantiles than the observed data at any site.”

And added the following to the main text:

“The model does not give higher quantiles than the observed data at any site.”

Have added a section “Suggestions for further work”, paraphrasing your comment in the following text:

“The model/observation departures seen in Fig. B1 have a systematic feature which is consistent over spatial regions, e.g., south-west and north-west UK. This suggests that it should be possible to account for these departures through a smooth spatial function which maps the differences in quantiles between the observations and model data. With this adjustment it is possible that the currently identified under-estimation may be corrected before making the tail-based GEV/GPD comparisons shown here.”

On first reading your suggestion, I thought: “that would only correct the location and scale parameters, which we are suggesting be taken from the observations anyway --- it wouldn’t make sense to correct the model shape parameters using the shorter observational records and then argue that the observational shape parameters should be replaced with those of the model”. But then I remembered about the scale-shape compensation which can occur at the fitting stage. Is that the reason for your suggestion?

Shape Parameters with Similar Spatial Patterns

One of your exciting findings is that the spatial pattern of the shape parameter estimates is similar for the CFB2018 estimates from observed data and your estimates from the model, but with a systematic bias between them. This suggests using an alternative to the CFB2018 approach by explicitly exploiting this finding. Let $\xi_{obs}(x)$ and $\xi_{model}(x)$ be the shape parameters for the observed and model data respectively, for all gauged sites x . What you are saying is that you believe in the spatial variation of $\xi_{model}(x)$ but not its mean value. So you believe that $\xi_{obs}(x) = \xi_{model}(x) + \xi_{bias}$, where ξ_{bias} is a fixed constant that does not vary over x . Your estimates indicate that $\xi_{bias} > 0$. Fixing estimates of $\xi_{model}(x)$ for all x and fitting the function of $\xi_{obs}(x)$ over x now means only one parameter, ξ_{bias} , needs estimating. This could lead to substantial reductions in the uncertainty of $\xi_{obs}(x)$ estimates over x .

I like that idea a lot! Have incorporated it in the main text. Hope that is OK. Have acknowledged your help in the acknowledgements section.

Line 344: have removed the phrase “without the subjectivity of a prior” and added the following:

“Given the need for some kind of constraint on the shape parameter when fitting observational records, use of shape parameters from a long simulation holds the promise of reducing uncertainties. For example, if we assume that the model-diagnosed spatial pattern of shape parameters is correct but uniformly biased by a scalar ξ_{bias} (which does not vary over sites), $\xi_{\text{true}}(x) = \xi_{\text{model}}(x) + \xi_{\text{bias}}$, where x is a vector of sites, then we can use the observations from *all sites* to estimate the one scalar parameter ξ_{bias} . This could lead to substantial reductions in the uncertainty of ξ_{true} estimates over x .”

Penalised Likelihood

The shape parameter estimates in Figure 2 (d) based on the 483 years of model data show some site-to-site variations which are more pronounced than the broader smooth variations across coastlines. This suggests that they would also benefit from the penalty-based approach used in CFB2018 work.

Thanks, have paraphrased this in the new “Suggestions for further work” section.

At a few points you state that the prior/penalty is subjective and that this is a disadvantage. Yet you also point that the smooth pattern of the shape parameter estimates this gives for the observed data agrees well with the similar unpenalised estimates using model data. You say this is a really positive feature of the model data, we would also take this as supporting the value of the process of creating penalised estimates from the observed data. The penalty is giving something meaningful, so the effect of the claimed “subjectivity” is positive for the observed data analysis.

I completely agree. I’m sorry if this did not come across clearly in the draft. Indeed, in some other experiments (not discussed in the paper) I tried unconstrained GEV fits to the annual maxima from the observations. This spoiled the strong correlation with the model-diagnosed shape parameters, so yes, completely agree that the CFB2018 penalisation process adds value over an unconstrained fit.

In the draft, we have the following text:

This strong correlation between the two spatial patterns of shape parameter diagnosed from independent sources (i.e. our model simulation and the tide-gauge data) is remarkable. It both supports the spatial pattern of the shape parameter as a real, physically-determined phenomenon (as opposed to a statistical artefact), **and gives further credibility to both the CFB2018 approach and our model.** [bold font is not in the draft].

The prior that was selected in the CFB2018 work was not subjective in the traditional sense of a subjective prior in Bayesian methods. It was actually a data-based prior which corresponds to an empirical Bayesian prior, using all the information that separately estimated shape parameters for UK skew surge provide. The effect of this was simply to move shape parameter estimates more towards the UK average, with the larger changes coming for sites with shorter record lengths.

Thank you for your guidance. I didn’t previously understand that difference in the definition of subjective vs empirical. Have added your description to the draft as follows, and removed all references to “subjectivity”

“CFB2018 employed a data-based prior using all the information that separately estimated shape parameters for UK skew surge provide. The effect of this was simply to move shape parameter estimates more towards the UK average, with the larger changes coming for sites with shorter record lengths.”

Investigation of Interaction between Tide and Skew Surge

We are really pleased to see the use of numerical models to explore systematically the widely claimed property that skew surges are independent of their associated peak tidal values. The work in Section 5 where you explore the timing effects of skew surge events relative to tides is very illuminating. It adds greatly to the existing empirical evidence for this feature at Sheerness. It would be interesting to have your opinions about what interactions are expected elsewhere in the UK given that the CFB2018 estimates over UK sites all assume independence. Since this analysis is based on simulations from HadGEM3-GC3-MM, this presents a physically-based justification for dependence between skew surge and peak tide at Sheerness. It suggests future research is required to investigate this. It may be that interaction occurs to some extent everywhere, but at a practical level it is not important apart from some locations.

Agreed. Have added a sentence to the “Suggestions for further work” section. For what it’s worth, my guess is that interaction is stronger at Sheerness than elsewhere, and is associated with the surge and tide travelling a long way together (all the way down the east coast) in shallow water. But that is just a guess.

In current work we have been investigating this feature empirically at a limited number of sites. Here we report some of the methods and findings using data from the tide gauges at Sheerness and Heysham (located close to Workington which you analyse). We study data from Heysham in 1964-2016, of which 17.5% is missing, and 1980-2016 at Sheerness, where 9.1% is missing.

We define extreme skew surges as exceedances of the 0.95 quantile at each site. Figure 1 shows scatter plots of extreme skew surges against their associated ranked peak tide; these are ranked so that 1 corresponds to the smallest observation and n is the largest, where n is the total number of observations. If the two components were independent, we expect extreme skew surges to be uniformly distributed over ranks. We test this using a Kolmogorov-Smirnov test for uniformity; the p value at Heysham is 0.0224 whilst at Sheerness it is 1.40×10^{-7} . Clearly the p value at Sheerness is much smaller, and provides statistical evidence that extreme skew surges are not independent of peak tide. However, at Heysham, there is not sufficient evidence to reject the claim of independence at the 0.01 significance level. This is clear in Figure 1, where more extreme skew surges occur on lower tides at Sheerness - this agrees with your findings and those of Williams et al. (2016).

We investigate this further by looking at skew surge and peak tide dependence on a month-by-month basis, using ideas from Williams et al. (2016). Figure 2 compares the distribution of peak tides associated

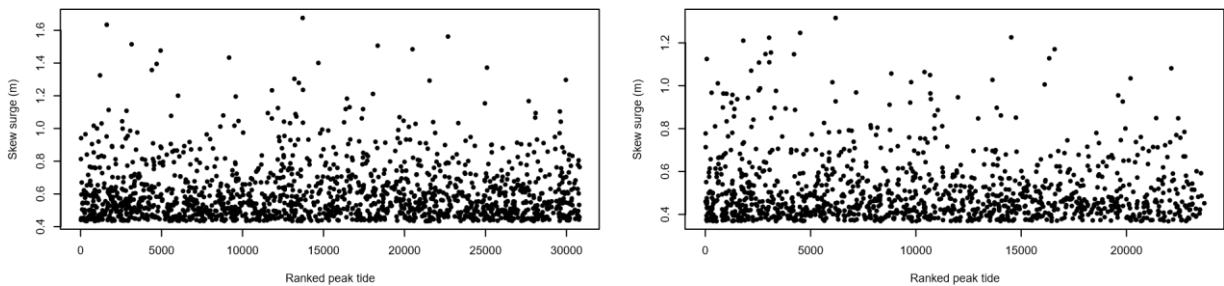


Figure 1: Extreme skew surge observations against ranked peak tide at Heysham (left) and Sheerness (right).

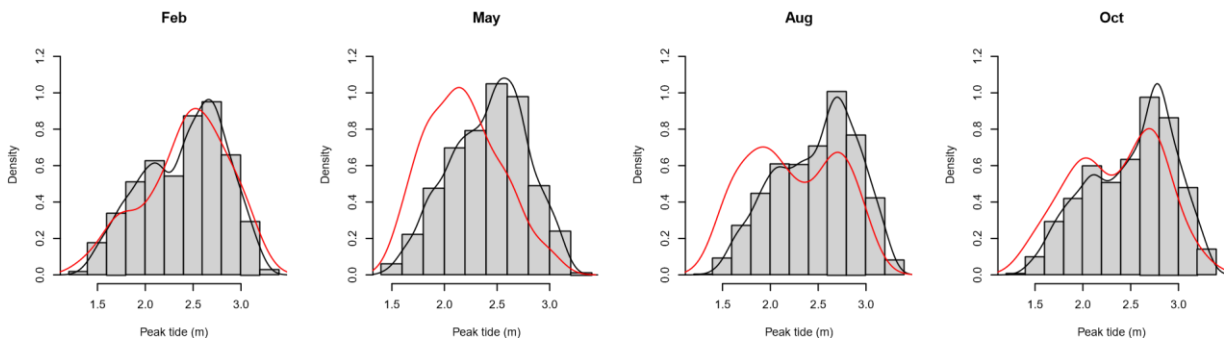


Figure 2: Monthly distributions of peak tides at Sheerness in February, May, August and October. The probability density function of all peak tides (black) and peak tides associated with extreme skew surge (red) are interpolated onto each distribution.

with all skew surges and the distribution of peak tides associated with extreme skew surges for February, May, August and October. If peak tide and skew surge are independent, these two distributions should be the same, up to sampling variation. We estimate the probability density function (pdf) using a Gaussian kernel density estimate, and use an Anderson-Darling test to check if peak tides come from the same distribution as peak tides associated with extreme skew surges. Figure 2 highlights how the dependence between skew surge and peak tide is changing with the time of year; the distribution of peak tides and the peak tides associated with extreme skew surge are most different in May and least different in February. In May, the mode of the distribution of peaks tides associated with extreme skew surges has shifted to a lower value than the distribution of all peak tides. Results from the Anderson Darling test for every month tell us there is insufficient evidence to reject the null hypothesis that peak tides and the peak tides associated with extreme skew surge come from the same distribution in February, March, September, November and December. In these months, we conclude it is reasonable to assume skew surge and peak tide are independent. In the remaining months, we find sufficient evidence to suggest the two components are dependent. We believe this poses a really interesting area for further research.

In D'Arcy et al. (2021) we fit a GPD (see Coles (2001) for details) to extreme skew surges that accounts for seasonal variations, since they are more extreme in the winter. Our results show this is an improvement on a standard GPD fit, where skew surges are assumed to be independent of peak tide. Here, we investigate adding a covariate of peak tide into our skew surge model to account for the dependence found at Sheerness. D'Arcy et al. (2021) defines extreme skew surges as exceedances of the monthly 0.95 quantile. Non-stationarity is accounted for in the scale parameter of the GPD through a daily covariate $d = 1, \dots, 365$:

$$\sigma_d = a + b \sin\left(\frac{2\pi}{365}(d - \phi)\right) \quad (3.1)$$

for $a, b, \phi \in \mathbb{R}$ parameters to be estimated. Note that the shape parameter of the GPD is fixed across months. To account for the dependence between skew surge and tide, we now also consider the following parameterisation on the scale parameter:

$$\sigma_{d,t} = a + b \sin\left(\frac{2\pi}{365}(d - \phi)\right) + ct \quad (3.2)$$

where $c \in \mathbb{R}$ is another parameter to be estimated, and t the associated peak tide observation. We fit a GPD with both (3.1) and (3.2) formulations to extreme skew surges at Heysham and Sheerness. The Akaike information criteria (AIC), frequently used for model selection, suggests that formulation (3.1) gives a better model fit at Heysham, i.e., an independence conclusion is supported by this analysis. Whereas, AIC suggests that the parametrisation in equation (3.2) yields a better fit at Sheerness. By ordering peak tide observations from smallest to largest and calculating $\sigma_{d,t}$ at its winter peak ($d = 365$) for each value, we observe a 15% reduction in the scale parameter as tides increase at Sheerness, which corresponds to an equal percentage reduction in the skew surge quantiles for excesses of the December skew surge threshold. We also compare the model fits at each site using a Likelihood Ratio Test. At Heysham the p value is 0.657 which provides insufficient evidence to reject the null hypothesis, that is the simpler model (in equation (3.1)) is sufficient. Whereas, at Sheerness we get a much smaller p value of 0.059, which provides statistically significant evidence at the 0.1 significance level to reject the null hypothesis and conclude that the more complex model is required. This highlights the importance of accounting for skew surge and peak tide dependence at sites where the independence assumption is not justified.

This is very interesting. Thank you for including it in your review. In the draft I have added a citation to your forthcoming publication.

Ungauged Sites

We feel you do not do full justice to your developments given the focus is on comparisons made at gauged sites with long and trustworthy records. The real value in modelled data is the ability to give estimates at other sites. This could be the focus of a natural follow up paper.

Have added a sentence in the “Suggestions for further work”

Other Comments

5. In the abstract, you say “results suggest an event of this magnitude has an expected frequency of about 1 in 500 years at [Sheerness]” when referring to the North Sea floods of 1953. In the paper, the only result to show this is presented in Section 6.2: “the fact that the 483-year surge-only simulation produces more than one event of comparable magnitude to the simulated 1953 event suggest the return period of the 1953 atmospheric forcing is less than 483 years.” If this is the justification, it would be better to more formally quantify this using your statistical models.

I can see several different ways to approach that. In the preceding sections we have advocated using only the shape parameter from the simulations. Using the observational location and scale parameters (as in Fig. 2 panels (a) and (b)) and the simulation-diagnosed shape parameter (panel (d)), the range of the Wadey skew surges as shown in our figure 7 correspond to return periods ranging between about 650 and 2000 years. (Note to self: code_ABR). Having said that, I feel that this level of detail is not appropriate in this preliminary “Towards using state-of-the-art climate models...” type of publication, so I have cut the reference to the expected frequency out of the abstract, and cut the corresponding short paragraph from the main text.

6. In the Introduction (line 51) you list assumptions required to fit extreme value models as ‘events are effectively random and statistically independent of each other.’ But what about events being identically distributed (or stationary), it is clear that tide and storm surge (or skew surge) are not stationary as both process exhibit seasonality. Instead of “effectively random” it would be more mathematically correct to say “stochastic.” Extreme value methods also handle dependence, so “independent of each other” is not formally required.

Rephrased as follows:

The statistical models which are fitted to the observational data in order to infer the levels of unprecedented extremes are supported by mathematical arguments which may require assumptions such as the assumption that the events are stochastic. We know that the real-world events are deterministic, and furthermore may be auto-correlated over a range of timescales. Such auto-correlation can be accounted for within the statistical framework, for example by the use of an extremal index~\citep{Tawn1992estimating, Batstone2013UK}. Alternatively, a physically-based numerical model has the potential to directly address both determinism and auto-correlation by simulating them.

7. On line 84 climate change is discussed. Tide gauge observations will exhibit an approximately linear mean trend due to sea level rise, it is unclear whether this has been removed before comparison with the simulations in Section 4. It is likely that removing this change will not change the results significantly, but it is important to remove this non-stationary effect before fitting extreme value models.

Have added this phrase:

“The trend due to sea level rise can be seen in the tide gauge observations and was carefully removed before making a statistical fit to the extremes. For details see CFB2018. Our numerical model of the shelf sea does not include any change in mean sea level.”

8. Figure C1 shows the reduction in uncertainty on shape parameter when record lengths are increased in the tide gauge network, but it would be nice to see this for more record lengths that are equally spaced

(say 10 to 500 years in increments of 10 years) to really highlight this - with measures of uncertainty as in Figure C1 (d).

Good idea. Done. Thanks.

Technical Corrections

10. Table 1 gives a list of useful acronyms and symbols, it would help to include what 'GC3' and 'MM' stand for in 'HadGEM3-GC3-MM' as it is not mentioned in text.

Thank you for pointing this out. Have added it to Table 1.

11. Equation (2) has a surplus close bracket. Corrected, thanks.

12. Equation (6), should this derivative be evaluated at $L = \log(u)$ rather than $y = u$ since y is the return level (as in equation (1))?

No, I don't think so. However I acknowledge that, in equation 6, the point on the RL curve at which to evaluate the gradient is specified in a non-standard way: in terms of the ordinate (y) instead of the usual specification in terms of the abscissa (L). Your suggested phrase "evaluated at $L = \log(u)$ " does not mean the same thing. Please see the uploaded supplement **equation6_revisited.pdf** for further explanation.

Note to self: upload eq 6 supplement.

13. Figures 2 (c) and (d) would benefit from having 95% confidence intervals on shape parameters using the model data.

Have added these, and commented in the text on the added value in terms of certainty that the CFB method provides compared to a simple GEV fit to annual maxima --- this is well-illustrated by the confidence intervals shown.

14. Line 283: We think the wrong figure is referenced here.

Please note this is figure E.1 in the CFB technical publication ("their figure E.1", as stated in the text). In the copy I downloaded from here:

https://assets.publishing.service.gov.uk/media/603652cce90e0740b7caac9d/Coastal_flood_boundary_conditions_for_the_UK_2018_update_-_technical_report.pdf

figure E.1 is on page 64, titled "Estimated shape parameter for the UK tide gauges."

15. Line 319: "They usedtime-series." This has nothing to do with how CFB2018 estimate the shape parameters and should be cut. It only links to the derivation of SWL return levels.

Have cut these two sentences as advised.

16. Line 349: You hint that the reason for the disagreement between the methods could be due the short observed data series. But Figure B1 shows major departures between the observed data and the model data in the body of the distributions to a level that could in no way be due to short observed series and must be that the model data does not reproduce well the observed data. Linked to this, in discussing Figure 1 you say the agreement at both sites is "excellent". With the amount of data at Sheerness it is clear that there are major disagreements that cannot be accounted by sampling variations.

Have reworded that paragraph to be less bullish:

"Figure 1 shows good model vs observations agreement at Workington, and even these two sites alone illustrate that our modelling system is able to simulate unprecedented skew surge events (i.e. events of a magnitude not found in the tide-gauge record). However, the quality of agreement shown at Workington

is not exhibited everywhere. Empirical return level plots of skew surge for a set of 44 tide gauge locations around the UK coastline are shown in the appendix in Fig. B1. This gives a qualitative, visual sense of the realism of the model in terms of the simulated extremes. The good agreement at Workington can be contrasted with the poor agreement at, for example, Newlyn or Aberdeen, where the simulated extremes are negatively biased relative to the corresponding observations. The model does not give higher quantiles than the observed data at any site. We argue (below) that the simulation may nevertheless be able to add value to estimations of unprecedented events, even where a bias exists.”

17. Line 421: You suddenly mention a “kernel” to spread the duration of events but fail to provide any information. A little detail would be helpful here.

I have added the following sentences explaining the choice and purpose of the kernel:

“The kernel was designed to represent the important features of the RACMO-driven surge-only simulation, i.e., the approximate duration and shape of the time series plot. The purpose of convolution with the kernel is to identify those events which correlate well (in terms of their time series plot) with the RACMO-driven simulation, in other words, events which not only produce a large surge, but are also of comparable duration to the RACMO-driven simulation. The kernel was not used to modify events, but simply to identify significant ones.”

18. Figure E2: A description of the different points would be useful.

Have added the following text: “The blue points show simulated skew surge against simulated astronomical high water for the 16 extreme atmospheric events with timing adjusted such that the event coincided with a simulated spring tide and a simulated neap tide, as described above. The artificial case of no tide is also shown. Grey points are as in Williams (2016).”

References

Coles, S. G. (2001). *An Introduction to Statistical Modeling of Extreme Values*. Springer, London.

D’Arcy, E., Tawn, J., Joly-Laugel, A., and Sifnioti, D. E. (2021). Accounting for seasonality in extreme sea level estimation (*in preparation*).

Williams, J., Horsburgh, K. J., Williams, J. A., and Proctor, R. N. (2016). Tide and skew surge independence: New insights for flood risk. *Geophysical Research Letters*, 43(12):6410–6417.

AC4:

This was an earlier draft of AC5, including a misleading typo. Have included AC5 instead here.

AC5:

Equation 6 revisited

Tom Howard

September 9, 2021

Intro

The purpose of this document is to give a more detailed explanation of Equation 6 of “Towards using state-of-the-art climate models to help constrain estimates of unprecedented UK storm surges”, Howard and Williams (2021).

Statistical Modelling of Extreme Values

To identify, for example, the 1000-year return level based solely on tide-gauge observations, some philosophy for making out-of-sample estimates is required. The usual approach is to exploit the most extreme observations, and theories concerning their behaviour, under some restrictive assumptions.

Annual Maxima

One popular and simple approach is fitting a Generalised Extreme Value (GEV) distribution to the annual maxima. The GEV distribution (GEVD) arises as the limiting case for block maxima as the block size tends to infinity. In the case of annual maxima, “block” means one year. The GEVD is characterised by three parameters. For readers unfamiliar with the GEVD, it may be helpful to picture the effect of these parameters in terms of a return-level curve, such as the ones shown in Fig. ???. The location parameter, μ , is comparable to an intercept. An increase in μ slides the whole curve up the Y-axis. μ is the Y-value (return level) evaluated at the one-year return period:

$$\mu = y \Big|_{L=0}$$

where $L = \log(\text{return period})$ and y is the return level. Notice that, though not particularly useful, this could be written

$$\mu = y \Big|_{y=\mu}$$

The GEV scale parameter, σ , is the gradient of the curve, evaluated at the one-year return period. This could either be written as

$$\sigma = \frac{dy}{dL} \Big|_{L=0} \quad (1)$$

or, for comparison with equation 6,

$$\sigma = \frac{dy}{dL} \Big|_{y=\mu}$$

Since y is a monotonic function of L , and in view of the first (unnumbered) equation, this is an alternative way to unambiguously define the point on the RL curve at which to evaluate the gradient. It’s just specified in a non-standard way: in terms of the ordinate (y) instead of the usual specification in terms of the abscissa (L).

The shape parameter, ξ , determines the curvature. Negative ξ corresponds to a curve which flattens out at high return periods, approaching an upper bound as the return period tends to infinity. With positive ξ the curve has no upper bound, but has a lower bound as the return level decreases. When $\xi = 0$ the curve is a straight line and has neither lower nor upper bound. This follows the convention of [?] for the shape parameter. However, not all sources follow this convention. In CFB2018, “shape parameter” refers to the negative of our ξ . In the wider literature the “shape parameter” may refer to the negative or the reciprocal of our ξ . To make our shape parameter notation unambiguous: if Y is a random variable with GEV distribution, our shape parameter ξ is defined such that the distribution of Y is given by

$$P(Y < y) = \exp \left\{ - \left[1 + \xi \left(\frac{y - \mu}{\sigma} \right) \right]^{-1/\xi} \right\} \quad (2)$$

This can be more simply expressed as the corresponding return level curve, which is

$$\frac{y - \mu}{\sigma} = \frac{R^\xi - 1}{\xi} \quad (3)$$

where the average recurrence interval (or “return period”) is R and the corresponding return level is y . The connection between equations 2 and 3 is seen by regarding exceedances of the R -year return level y as Poisson-distributed random occurrences, occurring at an average rate

$$\lambda = 1/R \quad (4)$$

The probability of no such occurrences in a given year is then given by standard Poisson statistics:

$$P(\text{no occurrences}) = P(Y < y) = \exp(-\lambda) \quad (5)$$

Combining 3, 4 and 5 gives equation 2. The particular case $\xi = 0$ is obtained by taking the limit as $\xi \rightarrow 0$.

Peaks over Threshold

The most extreme storm surges in the UK are caused by the storminess of the winter atmosphere, so the annual maximum event is always expected to occur in winter. Thus, an advantage of the annual-maxima approach described above is that the annual maxima are typically very well separated from each other and thus can be considered independent, particularly if the nominal year change is taken to be in the summer. A disadvantage of the approach is that it uses only the annual maxima. On the other hand, the peaks-over-threshold (POT) approach uses all of the data exceeding a chosen threshold. This formed part of the approach taken by CFB2018. An advantage of this approach is that, if a low-enough threshold is used, it has the potential to exploit more of the available data (i.e. an average of more than one extreme event per year), whilst including only extreme events. Such exploitation of more data usually reduces the uncertainties in inferred statistics (e.g. the out-of-sample estimates). This is particularly desirable when short observational records limit the available extremes. However, if the threshold is too low, some of the data included can no longer be considered “extreme” and may bias the result. This is the wellrecognised bias-variance trade-off. Another disadvantage is that including more than one event from a winter may compromise the independence of the events. (Skew surge can be evaluated for every high tide, and a weather system can generate a substantial skew surge on successive high tides.) Dependence is accommodated by CFB2018 using an extremal index... For a detailed comparison of the annual-maxima and POT approaches see...

The usual POT approach is to fit a Generalised Pareto Distribution (GPD) to the peaks. The GPD has two parameters. The shape parameter ξ is shared with the GEVD. The GPD scale parameter, σ , is the gradient of the plot of e

return level against log of return period at the return period of the chosen threshold, u ,

$$\tilde{\sigma} = \left. \frac{dy}{dL} \right|_{y=u} \quad (6)$$

As in the unnumbered equation following equation 1, the point on the RL curve at which to evaluate the gradient is specified in a non-standard way: in terms of the ordinate (y) instead of the usual specification in terms of the abscissa (L). σ is a property of both the extreme value distribution and the chosen threshold. The GEV scale parameter, σ , on the other hand, is a property of the extreme value distribution only and is thus a more fundamental parameter for making comparisons: it can be used in a like-for-like comparison of the results of different thresholds, or for comparison of GEV and GPD results. The two different scale parameters are related by $\sigma = \tilde{\sigma} \lambda_u^\xi$, where λ_u is the expected number of exceedances of u per year.

Though not formally a parameter of the GPD, a threshold must be chosen. CFB2018 tested 14 different thresholds and, finding no clear support for dismissal of any, elected to evaluate statistics based on each threshold and identify the median as the best estimate.

RC4:

Thank you for the positive way you have responded to our suggestions and also for the extra clarification on the return level curve connections with the scale parameter.

Eleanor D'Arcy and Jonathan Tawn

AC5: Thank you again for your review.

Latexdiff showing differences between original submission and revised submission following all review comments follows (assuming I manage to combine it successfully)...

Towards using state-of-the-art climate models to help constrain estimates of unprecedented UK storm surges

Tom Howard^{Met Office Hadley Centre} and Simon David Paul Williams^{National Oceanography Centre}

¹Met Office, FitzRoy Road, Exeter EX1 3PB, UK

²National Oceanography Centre, Joseph Proudman Building, 6 Brownlow Street, Liverpool L3 5DA, UK

Correspondence: Tom Howard (tom.howard@metoffice.gov.uk)

Abstract.

Our ability to quantify the likelihood of present-day extreme sea level (ESL) events is limited by the length of tide gauge records around the UK, and this results in substantial uncertainties in return level curves at many sites. In this work, we explore the potential for a state-of-the-art climate model, HadGEM3-GC3, to help refine our understanding of present-day coastal flood risk associated with extreme storm surges, which are the dominant driver of ESL events for the UK and wider European shelf seas.

We use a 483-year present-day control simulation from HadGEM3-GC3-MM (1/4 degree ocean, approx 60 km atmosphere in mid-latitudes) to drive a northwest European shelf seas model and generate a new dataset of simulated UK storm surges. The variable analysed is the skew surge (the difference between the high water level and the predicted astronomical high tide), which is widely used in analysis of storm surge events. The modelling system can simulate skew surge events comparable to the catastrophic 1953 North Sea storm surge, which resulted in widespread flooding, evacuation of 32 thousand people and hundreds of fatalities across the UK alone, along with many hundreds more in mainland Europe. Our model simulations show good agreement with an independent re-analysis of the 1953 surge event ~~and suggest that a skew surge event of this magnitude has an expected frequency of about 1 in 500 years~~ at the mouth of the river Thames. For that site, we also revisit the assumption of skew surge/tide independence. Our model results suggest that at that site for the most extreme surges, tide/surge interaction significantly attenuates extreme skew surges on a spring tide compared to a neap tide.

Around the UK coastline, the extreme tail shape parameters diagnosed from our simulation correlate very well (Pearson's r greater than 0.85), in terms of spatial variability, with those used in the UK government's current guidance (which are diagnosed from tide-gauge observations), but ours ~~can be diagnosed without the use of a subjective prior~~ have smaller uncertainties.

Despite the strong correlation, our diagnosed shape parameters are biased low relative to the current guidance. This bias is also seen when we replace HadGEM3-GC3-MM with a reanalysis, so we conclude that the bias is likely associated with limitations in the shelf sea model used here.

Overall, the work suggests that climate model simulations may prove useful as an additional line of evidence to inform assessments of present-day coastal flood risk.

25 *Copyright statement.* The works published in this journal are distributed under the Creative Commons Attribution 4.0 License. This licence does not affect the Crown copyright work, which is re-usable under the Open Government Licence (OGL). The Creative Commons Attribution 4.0 License and the OGL are interoperable and do not conflict with, reduce or limit each other. ©British Crown copyright 2021, the Met Office, UK. Also contains Environment Agency information ©Environment Agency and database right. This does not conflict with the CC BY 4.0 licence.

30 **1 Introduction**

Around £150 billion of assets and 4 million people in the UK are at risk from coastal flooding (Haigh et al., 2017), and estimated damages to the UK from coastal flooding are of the order of £500 million per year (Edwards, 2017). It is neither technically feasible nor economically affordable to prevent all such flooding, so policymakers use a risk-based approach. Typically, coastal flood protection is mandated based on an extreme high water “return level” with an estimated average
35 recurrence interval, which is the expected average time between exceedances of that level (conceived as averaged over a period including many such exceedances). The average recurrence interval is sometimes called a “return period”. We use these names interchangeably here, although some authors use a different definition of the return period. Typically, assets with high value and/or high vulnerability will have a mandate for protection against a return period of a thousand or even ten thousand years. Typical tide gauge records cover much shorter periods (of the order of 30 to 150 years). To address this, the traditional
40 approach is to fit a statistical extreme value model in order to extrapolate from the observations. Many different statistical approaches have been used (Haigh et al., 2010; Batstone et al., 2013); see §3.3. However, even using the current best practice, the inevitable extrapolation involved means that the uncertainties in the magnitude of very rare (perhaps unprecedented) events may be very large. For example, the size of the 90 % confidence interval on the 10,000-year return level at Sheerness is around 1.6 metres (Environment Agency, 2018). For comparison, the 90 % confidence interval on model projections of regional mean
45 sea level rise to 2100 relative to the 1981-2000 average under representative concentration pathway RCP8.5 for the same location is around 0.62 metres (Palmer et al., 2018). Coles (2001) discusses some of the advantages and disadvantages of the statistical modelling approach; he says: “*Caution is required in the interpretation of return level inferences especially for return levels corresponding to long return periods... estimates and their measures of precision are based on an assumption that the model is correct.*” The statistical models which are fitted to the observational data in order to infer the levels of unprecedented
50 extremes are supported by mathematical arguments which may require assumptions such as the assumption that the ~~events are effectively random, and statistically independent of each other~~ events are stochastic. We know that the real-world events are deterministic, and ~~may in reality furthermore may~~ be auto-correlated over a range of timescales. ~~Although some account can be taken of this~~ Such auto-correlation can be accounted for within the statistical framework, for example by the use of an extremal index (Tawn, 1992; Batstone et al., 2013), ~~a~~ Alternatively, a physically-based numerical model has the potential to
55 directly address ~~these issues~~ both determinism and auto-correlation by simulating them. Coles (2001) goes on to say: “*Though the [extreme value statistical] model is supported by mathematical argument, its use in extrapolation is based on unverifiable*

assumptions, and measures of uncertainty on return levels should properly be regarded as lower bounds that could be much greater if uncertainty due to model correctness were taken into account.”

60 An alternative approach is to exploit a physically-based numerical model of the coastal shelf waters. Such models typically parameterize the surface stress associated with winds and pressure from an atmospheric forecast model, and are routinely used to make short-range (e.g. less than 48-hour) forecasts of storm surges whenever a potentially hazardous atmospheric storm is identified in the atmospheric forecast. Bernier and Thompson (2006) found that when the atmospheric forecast model is replaced by an atmospheric hindcast, realistic extreme storm surge events were simulated in the northwest Atlantic.

65 Another approach is to make plausible modifications to the strength, track, or speed of selected observed atmospheric events, and use the resulting simulated atmospheric forcing to drive the coastal shelf model (Brown et al., 2010). Very recently Horsburgh et al. (2021) used this approach. They selected the storm of 5th December 2013 and made manual adjustments to the quasi-geostrophic potential vorticity field, inverting it to get dynamically self-consistent fields of sea-level pressure and wind. They showed that this approach can produce synthetic surges which are substantially larger than any in the observational record, for sites on the UK east coast. However this approach does not offer a way of quantifying the probability of the synthesized
70 events.

Yet another approach, adopted here and discussed further in §3.4, is to simulate extreme atmospheric events using a physically based numerical climate model, which in turn is used to drive the coastal shelf model.

An obvious advantage of this approach is that the model is based on verifiable real-world physics. Many climate model simulations extend over periods longer than the tide-gauge record. In particular, in order to evaluate model performance, 75 modellers use control simulations (with greenhouse gas forcing fixed at either pre-industrial or present-day levels) which may extend over many hundreds or even thousands of years. Ensemble simulations provide another potential source of data effectively covering a much longer period than the observations. Using the data from such simulations provides a further line of evidence in the effort to predict the magnitude and frequency of unprecedented events. [Van den Brink et al. \(2004\) used this approach to simulate storm surges at Hoek van Holland using the ECMWF seasonal forecast ensemble, successfully reducing the uncertainty in the 10000-year return level by a factor of four compared to using the observations alone.](#) This method was applied to seasonal rainfall totals in the UK (Thompson et al., 2017), and Grabemann et al. (2020) applied the method to extreme storm surges for locations in the German Bight, successfully identifying a number of simulated water levels exceeding those in the observational record since 1906.

85 This article reports a preliminary investigation into the value of using this approach to help form return level curves of storm surge around the UK coast with a view to providing improved likelihood information on the most extreme coastal water levels.

Climate change

Mean sea level is increasing, and will continue to increase, both at UK national scale (Palmer et al., 2018, 2020) and at global scale (Pörtner et al., 2019), and this will exacerbate future coastal flood risk. However, for many locations around the UK (exemplified by Sheerness as described above) the uncertainty in the projections of future mean sea-level rise is
90 not as large as the uncertainty associated with, say, the 1000-year return level of storm surge and it is the effort to reduce

this larger uncertainty that we are concerned with here: we are trying to “focus the snapshot” of conditions in the current climate. Thus we do not explicitly address mean sea level change in this work, but rather we note that the effects of mean sea level change and its uncertainty will need to be considered in addition to the present-day hazard which we discuss here, for example through the use of a sea-level rise allowance (Howard and Palmer, 2020). [The trend due to sea level rise can be seen in the tide gauge observations and was carefully removed before making a statistical fit to the extremes. For details see Environment Agency \(2018\). Our numerical model of the shelf sea does not include any change in mean sea level.](#) Also, we do not consider the effects of long-term change in the mean strength or location of the North Atlantic storm track (Shaw et al., 2016; Shepherd, 2014). Many studies (e.g. Palmer et al., 2018; Lowe et al., 2009; Sterl et al., 2009; Howard et al., 2019) have suggested that the change in local mean sea level will be the main contributor to the changes in the sea level extremes, as it has been in the past (Menéndez and Woodworth, 2010), with the change in the storm track making a smaller secondary contribution. We do not consider this secondary contribution here.

2 Nomenclature and Notation

For ease of reference, some terms which arise throughout this article are given in table 1.

3 Models, Methods and Data Sources

3.1 The CS3 coastal shelf model

Our barotropic coastal shelf model, CS3 (Continental Shelf 3, Horsburgh et al., 2008; Flather, 2000, 1994) is very similar to the CS3X (Continental Shelf 3 Extended) model which until very recently was used in the UK operational storm surge forecast/warning system. The domain and grid of CS3 are shown in the appendix in Fig. A1(a). The model produces a numerical solution of the discretized nonlinear shallow water equations with friction. The model is barotropic in the sense of solving the depth-averaged equations (i.e. it is a two-dimensional model). The horizontal resolution is approximately 1/9 degree latitude by 1/6 degree longitude (approximately 12 km). The model has been shown to perform particularly well during extreme storm surges in the southern North Sea (Horsburgh et al., 2008), forecasting surge in the Thames estuary to within 10 cm when driven by re-analysed meteorology. CS3 is “one of the most validated operational storm surge forecasting models in the world” (Horsburgh et al., 2021). Further details of storm surge model evaluation can be found in Furner et al. (2016); O’Neill et al. (2016); Palmer et al. (2018); Flather (2000). Typical RMS errors when forced with numerical weather prediction model atmospheric data are of the order of 10 cm.

3.2 Coastal Flood Boundary Conditions for the UK: update 2018

Coastal Flood Boundary Conditions for the UK: update 2018 (Environment Agency, 2018) (henceforth CFB2018) contains the latest UK government best estimates and uncertainty estimates for the distribution of extreme still water level (SWL) under present-day mean sea level. SWL can be thought of as the water level averaged over about five minutes to remove the short-

Acronym or Symbol	Description
CS3	Continental Shelf 3: our North-West European storm surge model. See §3.1.
HadGEMHadGEM3-GC3-MM	Hadley Centre Global Environment Model - See §3.4. in the Global Coupled configuration 3, Medium resolution atmosphere, Medium resolution ocean.
CFB2018	Coastal Flood Boundary Conditions for the UK: update 2018 (Environment Agency, 2018).
CMIP5	Climate Model Intercomparison Project, Phase 5 (Taylor et al., 2012).
CMIP6	Climate Model Intercomparison Project, Phase 6 (Eyring et al., 2016).
SWL	Still Water Level. Still water level includes the astronomical tides and surge but does not include the short-period oscillations due to waves. See §3.2.
GEV, GEVD	Generalised Extreme Value (Distribution) (Coles, 2001). See §3.3. Under appropriate conditions, annual maxima are expected to follow a GEVD.
GPD	Generalised Pareto Distribution (Coles, 2001). See §3.3. Under appropriate conditions, all extreme values over a high threshold are expected to follow a GPD.
POT	Peaks Over Threshold. A POT model uses all values over a high threshold (see GPD).
MLE	Maximum Likelihood Estimator (Coles, 2001). See §3.3.
PMLE; GMLE	Penalised Maximum Likelihood Estimator; Generalised Maximum Likelihood Estimator. Used synonymously here. See §3.3.
prior; penalty; constraint	These all refer to PMLE. See §3.3.
μ	GEV location parameter. See §3.3.
σ	GEV scale parameter. See §3.3.
ξ	Shape parameter. See §3.3.
$\tilde{\sigma}$	GPD scale parameter. See §3.3.
R	Return Period in years. See §3.3.
L	Natural logarithm of Return Period. See §3.3.
y	Return Level. See §3.3.

Table 1. Acronym and Symbols

period oscillations due to surface waves. It includes the astronomical tide and storm surge, and is the level that is reported at the tide gauges. The CFB2018 approach is based on data from tide-gauge observations, without reference to model simulations. Discussion of their approach is included below.

3.3 Statistical Modelling of Extreme Values

125 To identify, for example, the 1000-year return level based solely on tide-gauge observations, some philosophy for making out-of-sample estimates is required. The usual approach is to exploit the most extreme observations, and theories concerning their behaviour, under some restrictive assumptions.

Annual Maxima

One popular and simple approach is fitting a Generalised Extreme Value (GEV) distribution to the annual maxima. The GEV distribution (GEVD) arises as the limiting case for block maxima as the block size tends to infinity. In the case of annual maxima, “block” means one year. The GEVD is characterised by three parameters. For readers unfamiliar with the GEVD, it may be helpful to picture the effect of these parameters in terms of a return-level curve, such as the ones shown in Fig. 1. The location parameter, μ , is comparable to an intercept. An increase in μ slides the whole curve up the Y-axis. μ is the Y-value (return level) evaluated at the one-year return period. The GEV scale parameter, σ , is the gradient of the curve, evaluated at the one-year return period:

$$\sigma = \left. \frac{dy}{dL} \right|_{L=0} \quad (1)$$

where $L = \log(\text{return period})$ and y is the return level. The shape parameter, ξ , determines the curvature. Negative ξ corresponds to a curve which flattens out at high return periods, approaching an upper bound as the return period tends to infinity. With positive ξ the curve has no upper bound, but has a lower bound as the return level decreases. When $\xi = 0$ the curve is a straight line and has neither lower nor upper bound. This follows the convention of Coles (2001) for the shape parameter. However, not all sources follow this convention. In CFB2018, “shape parameter” refers to the negative of our ξ . In the wider literature the “shape parameter” may refer to the negative or the reciprocal of our ξ . To make our shape parameter notation unambiguous: if Y is a random variable with GEV distribution, our shape parameter ξ is defined such that the distribution of Y is given by

$$P(Y < y) = \exp \left\{ - \left[1 + \xi \left(\frac{y - \mu}{\sigma} \frac{y - \mu}{\sigma} \right) \right]^{-1/\xi} \right\} \quad (2)$$

This can be more simply expressed as the corresponding return level curve, which is

$$\frac{y - \mu}{\sigma} = \frac{R^\xi - 1}{\xi} \quad (3)$$

where the average recurrence interval (or “return period”) is R and the corresponding return level is y . The connection between equations 2 and 3 is seen by regarding exceedances of the R -year return level y as Poisson-distributed random occurrences, occurring at an average rate

$$\lambda = 1/R \quad (4)$$

The probability of no such occurrences in a given year is then given by standard Poisson statistics:

$$P(\text{no occurrences}) = P(Y < y) = \exp(-\lambda) \quad (5)$$

Combining 3, 4 and 5 gives equation 2. The particular case $\xi = 0$ is obtained by taking the limit as $\xi \rightarrow 0$.

The most extreme storm surges in the UK are caused by the storminess of the winter atmosphere, so the annual maximum event is always expected to occur in winter. Thus, an advantage of the annual-maxima approach described above is that the annual maxima are typically very well separated from each other and thus can be considered independent, particularly if the nominal year change is taken to be in the summer. A disadvantage of the approach is that it uses only the annual maxima.

160 On the other hand, the peaks-over-threshold (POT) approach uses all of the data exceeding a chosen threshold. This formed part of the approach taken by CFB2018 (Environment Agency, 2018). An advantage of this approach is that, if a low-enough threshold is used, it has the potential to exploit more of the available data (i.e. an average of more than one extreme event per year), whilst including only extreme events. Such exploitation of more data usually reduces the uncertainties in inferred statistics (e.g. the out-of-sample estimates). This is particularly desirable when short observational records limit the available

165 extremes. However, if the threshold is too low, some of the data included can no longer be considered “extreme” and may bias the result. This is the well-recognised bias-variance trade-off. Another disadvantage is that including more than one event from a winter may compromise the independence of the events. (Skew surge can be evaluated for every high tide, and a weather system can generate a substantial skew surge on successive high tides.) Dependence is accommodated by CFB2018 using an extremal index (Tawn, 1992; Batstone et al., 2013). For a detailed comparison of the annual-maxima and POT approaches

170 see Arns et al. (2013).

The usual POT approach is to fit a Generalised Pareto Distribution (GPD) to the peaks. The GPD has two parameters. The shape parameter ξ is shared with the GEVD. The GPD scale parameter, $\tilde{\sigma}$, is the gradient of the plot of return level against log of return period at the return period of the chosen threshold, u ,

$$\tilde{\sigma} = \left. \frac{dy}{dL} \right|_{y=u} \quad (6)$$

175 This is a property of both the extreme value distribution and the chosen threshold. The GEV scale parameter, σ , on the other hand, is a property of the extreme value distribution only and is thus a more fundamental parameter for making comparisons: it can be used in a like-for-like comparison of the results of different thresholds, or for comparison of GEV and GPD results. The two different scale parameters are related by $\sigma = \tilde{\sigma} \lambda_u^\xi$, where λ_u is the expected number of exceedances of u per year.

Though not formally a parameter of the GPD, a threshold must be chosen. CFB2018 tested 14 different thresholds and,

180 finding no clear support for dismissal of any, elected to evaluate statistics based on each threshold and identify the median as the best estimate.

Maximum Likelihood Estimation

As a model-fitting approach, CFB2018 adopted maximum likelihood estimation (MLE, Coles, 2001) and so do we.

Penalised Maximum Likelihood Estimation/Generalised Maximum Likelihood Estimation

185 A recognised problem of short records such as the relatively short tide-gauge record at some sites is the diagnosis of “noisy” and implausible shape parameters by MLE (see appendix C). We also show in appendix C that the uncertainty in estimating unprecedented events from observational records using MLE is dominated by uncertainty in the shape parameter. One fix for this is to put a ~~subjectively-chosen~~ prior (or “penalty function”) on the shape parameter (Coles and Dixon, 1999; Martins and Stedinger, 2000). This method was used by CFB2018. It is variously known as Generalised Maximum Likelihood Estimation
190 or Penalised Maximum Likelihood Estimation (PMLE). We also refer to the penalty function as a constraint. ~~We show below that the need for a penalty function is obviated in the case of our simulation, due to the long~~ CFB2018 employed a data-based prior using all the information that separately-estimated shape parameters for UK skew surge provide. The effect of this was simply to move shape parameter estimates more towards the UK average, with the larger changes coming for sites with shorter record lengths. ~~We argue that this removes some of the subjectivity.~~

195 Skew Surge Joint Probability Method

The large return-level uncertainties for long return-period events are mitigated by the use of the skew surge joint probability method, the current state-of-the-art approach. Extreme SWLs are composed of a high astronomical tide and a meteorological surge. The metric of choice for the meteorological component is the skew surge (de Vries et al., 1995): the difference between the (deterministic, predictable) astronomical high tide and the actual high water level (which typically arrives at a slightly
200 different time). See Palmer et al. (2018) for a schematic diagram illustrating the definition of skew surge (their figure A1.3.4). Under the assumption of tide-skew surge independence, which has substantial observational support (Williams et al., 2016), the level of the high tide is assumed to have no effect on the magnitude of the skew surge and thus any skew surge can combine with any high tide. This suggests a method (exploited by CFB2018) whereby the observed surges are decomposed into tide and skew surge to give a skew surge distribution, which can be convolved with the full, known distribution of high tides to
205 form the full distribution of high water levels. The extreme value modelling is only involved in establishing the high tail (i.e. the outside-sample part) of the skew surge distribution. The implication of this convolution is that although a very rare high water level might be a combination of an equally rare skew surge and an ordinary tide, it could also be formed by a very rare high tide (the distribution of which is well known) and an ordinary skew surge. A consequence is that uncertainties in very rare high water levels map to the uncertainties in less-rare skew surges, and these uncertainties are smaller than the uncertainties in
210 very rare skew surges. In other words, for a given return period, the high-water-level uncertainty is smaller than the skew-surge uncertainty. This is good, because it is the high water level that we are concerned about from a coastal flooding point of view. Having said all that, we do not have cause to use the skew surge joint probability method in this work; we only mention it due to its relevance to the CFB2018 estimates. We revisit the assumption of tide-skew surge independence in §5.2.

3.4 A free-running climate model as a driver of synthetic storm surges

215 The atmospheric jet over the north Atlantic, which is associated with the extratropical cyclones which drive surges on the
UK coast, has complex variability with a trimodal latitudinal behaviour (Woollings et al., 2010). A lot of effort in the climate
modelling community is directed to understanding and improving the quality of models' simulation of this behaviour, owing to
its importance in projections of climate change in the mid-latitudes (Shaw et al., 2016; Shepherd, 2014). Ongoing improvements
in the representation of the North Atlantic storm track in global climate models are discussed by Roberts et al. (2018); Priestley
220 et al. (2020).

Williams et al. (2015) show improvements in the representation of storm tracks in the CMIP6 (Eyring et al., 2016) gen-
eration Hadley Centre models relative to HadGEM2-AO (the Hadley Centre model which contributed to CMIP5), with both
HadGEM3-GC2 and HadGEM3-GC3 simulating the winter latitudinal variability well. Both models employ the ENDGame
revision to the dynamical core, which reduces the numerical damping associated with the semi-implicit advection scheme
225 and has been shown to increase synoptic variability (Williams et al., 2015). This suggests that a surge simulation driven by
HadGEM3-GC3 surface wind and pressure might yield realistic storm surges for the UK. We exploited a 483-year control
simulation of HadGEM3-GC3-MM (Williams et al., 2018, and references therein). In this simulation, greenhouse gas concen-
trations are fixed at pre-industrial levels. Its atmospheric horizontal resolution is N216 (approximately 60 km in mid-latitudes),
and ocean horizontal resolution approximately 1/4 degree. The atmospheric component (Walters et al., 2019) of this model
230 exhibits a very good representation of the storm track (as measured against the ERA Interim (Dee et al., 2011) reanalysis)
when forced with present-day sea-surface temperatures (pers. comm: Julia Lockwood, by email).

One argument that might be made against this approach is that the spatial resolution of the global climate model may be
inadequate to resolve all of the physical processes that might be important in generating extreme events, particularly small-
scale extremes. For example, contemporary global climate models do not have adequate resolution to synthesize a small
235 convective event such as a thunderstorm. However, three factors argue against this being a problem in the case of UK storm
surge modelling:

- Storm surge in the UK is usually driven by atmospheric baroclinic instability, which is a large-scale process, much larger
than the scale of a single thunderstorm, and well-captured by atmospheric models.
- Storm surge effectively integrates the driving atmospheric wind and pressure over a large area and time (Sterl et al.,
240 2009), so that shortcomings in the simulation of small temporal- or spatial-scale phenomena are relatively less important
than they would be in the simulation of, for example, localised short-duration extreme rainfall events.
- Storm surge generation occurs over the sea. It has been well-recognised for some time that the orographic drag schemes
used in atmospheric modelling improve the column-average windspeed at the expense of realistic surface windspeeds
over high ground (e.g. Howard and Clark, 2007). On account of this we might expect to find issues with simulation
245 of the surface wind climatology over land, particularly over high ground. Whilst this issue might also affect the ocean
points nearest to the coast, it has little effect over the open sea, where most of the surge is generated. de Winter et al.

(2013) evaluated 12 global climate models in terms of their simulation of wind over the North Sea, including two models from the HadGEM family. Their results show that these two models (along with two models from the GFDL-ESM family) exhibit a particularly realistic distribution of extreme winds (evaluated against a reanalysis), being well within the uncertainties of the reanalysis.

The UK Climate Projections 2018 Marine Report (Palmer et al., 2018) provides extensive evidence of the realism of storm surges simulated by CS3 when driven by climate model winds and pressure.

4 Results and Discussion

Example empirical return level plots of skew surge, comparing model and observational (tide-gauge) data at two sites are shown in Fig. 1. We take the model grid cell closest to the location of the real-world tide-gauge to represent that tide-gauge. The two sites chosen for this illustration are Sheerness, at the mouth of the river Thames in south-east England, a site of great economic and societal importance, and Workington, a coastal site in the north west of England which is typically affected by different storms (Haigh et al., 2016).



Figure 1. Empirical return level plots of skew surge, comparing simulated and observational (tide-gauge) data at Sheerness and Workington.

Figure 1 shows excellent model vs observations agreement for the two sites illustrated at Workington, and even these two sites alone illustrate that our modelling system is able to simulate unprecedented skew surge events (i.e. events of a magnitude not found in the tide-gauge record). However, the quality of agreement shown in Fig. 1 at Workington is not exhibited everywhere. Empirical return level plots of skew surge for a set of 44 tide gauge locations around the UK coastline are shown in the appendix in Fig. B1. This gives a qualitative, visual sense of the realism of the model in terms of the simulated extremes.

The ~~excellent agreement at Sheerness and~~ good agreement at Workington can be contrasted with the poor agreement at, for example, Newlyn or Aberdeen, where the simulated extremes are negatively biased relative to the corresponding observations. ~~We show below~~ The model does not give higher quantiles than the observed data at any site. We argue (below) that the simulation may nevertheless be able to add value to estimations of unprecedented events, even where a bias exists.

4.1 Quantitative evaluation of simulation of extremes

To make some quantification of the realism of the simulated extremes, we used the statistical models described in §3.3 to fit the simulated extremes. We fitted a GEV model (§3.3) to the simulated annual maxima using the MLE method (§3.3). We fitted the model pointwise (that is, for each tide gauge we fitted independently at a model grid cell closest to the tide-gauge of that site; this model grid cell is taken to represent that site). This gives a spatial distribution of diagnosed parameters. We find ~~excellent~~ good agreement between the simulation-based and tide-gauge-based location and scale parameters (Fig. 2). We also find a surprisingly good correlation between the spatial distribution of simulation-diagnosed shape parameters and the corresponding spatial distribution of shape parameters diagnosed by CFB2018. Pearson's r for the shape parameter correlation is 0.72 when we use a GEV fit to the simulated annual maxima, and 0.86 when we use a GPD fit to the simulated peaks over a threshold (see §4.2).

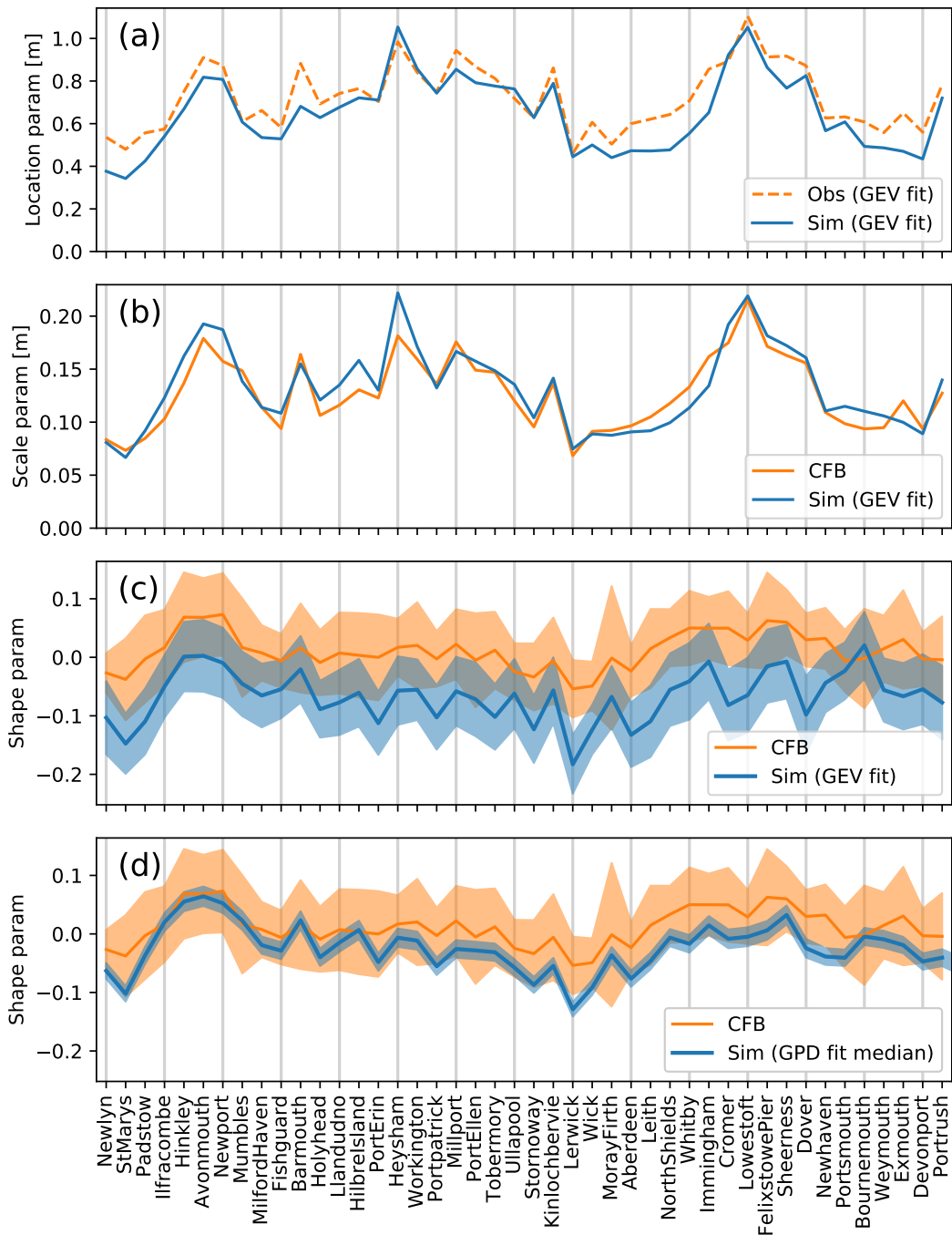


Figure 2. Comparison of simulation-based and observation-based skew surge extreme value distribution parameters. (a): Location parameter. (b) GEV scale parameter (σ). (c, d): Shape parameter. The correlation seen in all panels shows that the model successfully simulates the observed spatial variations in the extremes. [Pointwise uncertainties in the estimated shape parameters are included in panels \(c\) and \(d\). The CFB uncertainty shown is the 95% confidence interval of the GPD fit at a 95% threshold. In panel \(d\) \(the last panel\), the simulation uncertainty shown is evaluated in the same way. In panel \(c\) \(next-to-last\), the simulation uncertainty shown is the 95% confidence interval of the GEV fit to the simulated annual maxima.](#) For [further](#) details see main text.

A detailed description of Fig. 2 panels (a,b,c) follows. A complete description of panel (d) is deferred to §4.2. The correlation in all panels shows that the model successfully simulates the observed spatial variations in the skew surge extremes. In particular, good representation of the scale parameter at each site is important because this means that the temporal variability is well-simulated at that site. The absolute size of the scale parameter is significant, so we include zero in the Y-axis of panel (b). A scale parameter of zero would indicate no inter-annual variation in the extremes.

If we were to base our assessment on, for example, SWL relative to local Chart Datum (instead of skew surge), then the absolute value of the location parameter would have no particular significance: it would depend on a local offset. However, for skew surge the absolute value of the location parameter does have a significance: it represents a hypothetical absence of any atmospheric effects. For that reason we also include zero in the Y-axis of panel (a). A location parameter of zero would indicate no atmospheric effect on sea levels.

For all sites shown in Fig. 2, we obtained sufficient information regarding the CFB fit to the observations (pers. comm: Jenny Sansom, by email) to enable us to evaluate their GEV scale parameters (panel (b)). Their shape parameters (panel (c) and (d)) can be read from their figure E.1. For the location parameter (panel (a)) we used our own more crude fit to the tide-gauge annual maxima. We confirmed this crude fit against additional CFB information at a sample of nine sites. The crude fit was only used to estimate the observational location parameter.

Consideration of Fig. 2(a) shows that the simulation-diagnosed location parameters are in general slightly low compared to our crude estimate from the tide gauge data. At some sites this may be associated with locally poor representation of the details of the coast and bathymetry around the tide gauge due to the surge model resolution. However, the scale parameter (panel (b)) is generally in very good agreement with the CFB2018 results. This is reassuring because it indicates that the simulation is doing a good job of capturing the variability in the extremes (scale parameter), even though it shows an overall bias in the extremes (location parameter).

Panel (c) has ~~three main~~ the following features:

1. The shape parameters diagnosed from the simulation are well correlated with the CFB2018 shape parameters. This strong correlation between the two spatial patterns of shape parameter diagnosed from independent sources (i.e. our model simulation and the tide-gauge data) is remarkable. It both supports the spatial pattern of the shape parameter as a real, physically-determined phenomenon (as opposed to a statistical artefact), and gives further credibility to both the CFB2018 approach and our model. The authors are not aware of any previous work in which the spatial pattern of skew-surge shape parameter diagnosed from a simulation based on a free-running climate model has been shown to correlate well with the corresponding pattern diagnosed from observations.
2. The spread ~~of the~~ (i.e. the size of the spatial variations) of the shape parameters diagnosed by MLE fit (i.e. without constraint) to annual maxima from the simulation is comparable to ~~the spread that~~ of the CFB2018 shape parameters diagnosed by PMLE (i.e. with constraint), which in turn is similar to the spread of the ~~subjective~~ prior used by CFB2018. This again suggests that a long climate model simulation may be useful in constraining the shape parameters.

3. The pointwise shape parameter uncertainty (i.e. the uncertainty in the shape parameter at a given location) of the GEV fit to the simulated annual maxima is comparable to that of the constrained GPD fit to the observed surges above the 95% threshold, in spite of the shorter observational record lengths. This illustrates the added certainty of the CFB method over a simple GEV fit.
- 315 4. The fitted shape parameters for the simulation are more negative than the CFB2018 shape parameters. We return to this in §4.2.

The sites in Fig. 2 follow a clockwise orbit of the UK mainland coast starting at Newlyn in the south west, with the addition of St Mary's (Isles of Scilly), Port Erin (Isle of Man), Stornoway (Outer Hebrides), Lerwick (Shetland Isles) and Portrush (Northern Ireland). The sites are shown in Fig. A1(b) in the appendix.

320 4.2 Shape parameter

We return now to the fitted shape parameters for the simulation, which are more negative than the CFB2018 shape parameters. This is important because uncertainty in estimating unprecedented events from observational records using MLE is dominated by uncertainty in the shape parameter (see appendix C). This suggests that the shape parameter is the aspect where model simulations, with their long record lengths, may be able to help. We studied the negativity in several ways, with results which
325 are shown in Fig. 3 and Fig. 2 panels (c) and (d). Simulation results shown in Fig. 2 panels (a to c) are from a GEV fit to the simulated annual maxima, whereas CFB2018 results are from a GPD fit to the observations. To eliminate this potential source of difference, we also made a GPD fit to the simulation (Fig. (2) panel (d) and Fig. 3 line (C)). We applied the same treatment to simulations as was used by CFB2018 in order to make a like-for-like comparison. We did not apply a prior to produce the simulation shape parameters shown in Fig. 2 and Fig. 3 line (C).

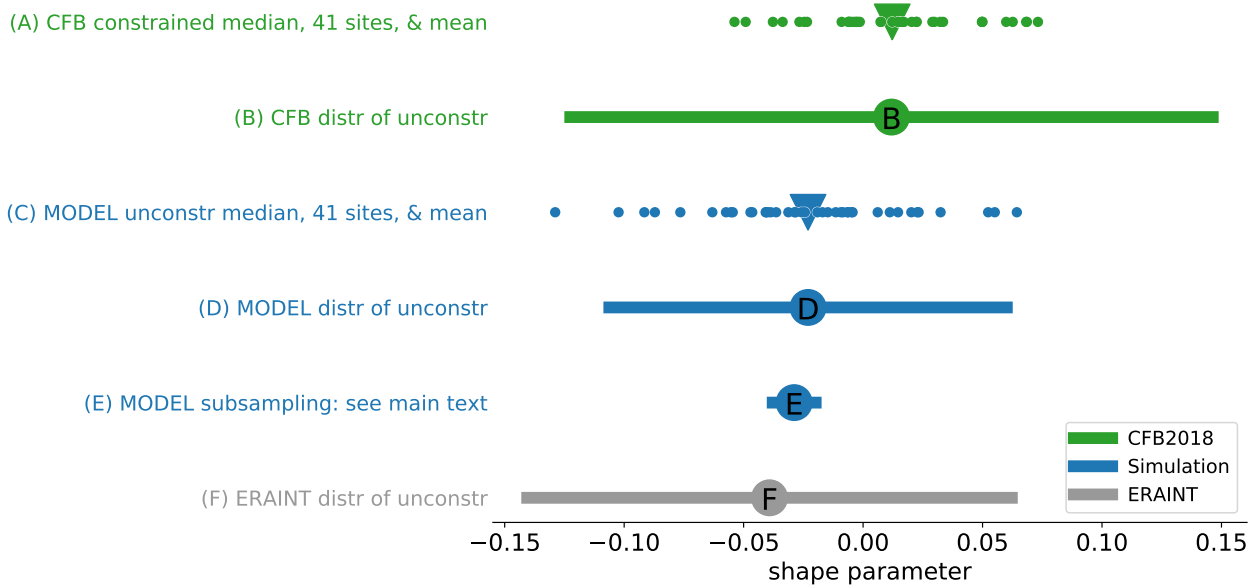


Figure 3. Further results related to the model vs CFB2018 shape parameter difference. Each line shows (X-axis; dimensionless) a distribution of GPD shape parameters, or a derived quantity such as the mean of several shape parameters. “unconstr” = Unconstrained. For full details see the main text.

330 Details of Fig. 3 follow. Lines (A) and (B) show shape parameter results from CFB2018. ~~They used a GPD fit to skew surges exceeding a threshold. They used an extremal index (Tawn, 1992; Batstone et al., 2013) to accommodate dependence in the time-series.~~ To compensate for the short observational record lengths, they used a prior to constrain the shape parameter (see §3.3). The prior (or “penalty function”) in turn was chosen by expert judgement informed by unconstrained GPD fits to the tide-gauge data. 14 different thresholds were tested, and results for each site and each threshold were pooled to form a distribution

335 of unconstrained shape parameters. This distribution is shown in line (B). The green line shows the range (characterised by two standard deviations either side of the mean). The filled disc labelled "B" shows the mean. The CFB2018 prior was taken to be a normal with the same mean but half the standard deviation. For each site, the finalised (constrained) shape parameter diagnosed by CFB2018 was chosen as the median of the PMLE results of the 14 different thresholds. These finalised shape parameters, one for each of 41 sites, are shown by the dots in line (A). (This same information is contained in Fig. 2 panels (c) and (d)). The

340 mean of these 41 shapes is shown by the filled triangle. We applied a similar approach (but without the prior) to our simulated skew surges to give the results shown in line (C). The Pearson’s r correlation between the model-diagnosed shape parameters using GPD fits (line C) and the CFB2018 shape parameters (line A) is 0.86. Owing to the much greater record length of the simulation, we did not need to apply any constraint to obtain the data on line (C).

345 The data of line (C) vs the data of line (A) show our most like-for-like simulation-vs-CFB2018 shape parameter comparison, and are also presented in Fig. 2 panel (d). That panel also shows our estimate of the uncertainty in the CFB2018 shape parameters, expressed as a 95% confidence interval. This estimate is based on the standard error of the CFB2018 fit at the 95% threshold (pers. comm: Jenny Sansom, by email). It is not straightforward to estimate the uncertainty of the CFB2018 shape parameters owing to their use of a median over results based on different thresholds ranging from 90% to 99%, but we suggest that the uncertainty in their 95% threshold result is representative.

350 Line (D) represents the full distribution of shape parameters diagnosed by GPD fit to the 483-year HadGEM3-driven model simulation. The blue line shows the range (characterised by two standard deviations either side of the mean). The range (as in line B) comes from variations in site and threshold used. The filled disc labelled “D” shows the mean. Thus D (model) corresponds to B (tide gauge). Clearly the 483-year model-diagnosed shapes are more negative than those derived from the shorter tide gauge data.

355 Given the need for some kind of constraint on the shape parameter when fitting observational records, use of shape parameters from a long simulation holds the promise of reducing uncertainties ~~without the subjectivity of a prior. Thus, the~~ For example, if we assume that the model-diagnosed spatial pattern of shape parameters is correct but uniformly biased by a scalar ξ_{bias} (which does not vary over sites), $\xi_{true}(x) = \xi_{model}(x) + \xi_{bias}$ where x is a vector of sites, then we can use the observations from all sites to estimate the one scalar parameter ξ_{bias} . This could lead to substantial reductions in the uncertainty of ξ_{true} estimates over x .

360 The more-negative shape parameters diagnosed by fitting the model data are, potentially, our most important finding, but further work is required to better understand the causes of this negativity. On one hand, it could be that limitations in the realism of either the atmospheric or the coastal shelf modelling distort the distributional tail relative to the real world. On the other hand, it could be that the physically-based model simulation gives better guidance on the distributional tail of the atmospheric storms which drive surges than does a statistical fit to the relatively short observational record of the surges themselves. In favour of the simulation, we can say that the emergence of realistic long-period natural variability in climate model simulations suggests their suitability for generating samples outside the observational record length. If it could be shown that the long-period variability in the simulation envelopes the observational results, this would give much stronger support to the use of the simulation.

370 Could it be, then, that if the simulation were sub-sampled in shorter periods to match the tide gauge record lengths, the value of a new metric (call it D' , the mean of the distribution of shape parameters diagnosed by GPD fit to the shorter sub-sampled HadGEM3-driven model simulation) would vary substantially so as to sometimes include values as large as “B”? To answer this question we sub-sampled the model many times to give a distribution of D' . This distribution is represented by line (E): the blue line shows the range of values of D' (characterised by two standard deviations either side of the mean; this range comes from random variations which we have introduced into the start time of the sub-samples, so that each sub-sample represents a randomly-chosen different "era" of the simulation) and the filled blue disc labelled “E” shows the mean value of D' . It is clear that this distribution *does not* include values as large as “B”, meaning that the *apparent positive* shape parameter bias of

the tide gauge results (B) relative to the model results (D) is *not* simply a “sampling error” associated with the shorter record lengths of the tide gauges, but rather a *real negative* bias in the model shape parameters relative to the tide gauges.

380 However, line (F) shows the distribution of unconstrained shape parameters diagnosed from a 29-year CS3 run forced by atmospheric surface wind and pressure based on the ERA-interim atmospheric reanalysis (Dee et al., 2011) that has been downscaled with the Swedish Meteorological and Hydrological Institute (SMHI) Rossby Centre regional atmospheric model (RCA4) as part of the Euro-CORDEX experiment (Jacob et al., 2014). This distribution shows at least as much negative bias as the HadGEM3-driven model simulation, even though the ERA reanalyses are widely viewed as the gold standard in terms
385 of representing the storminess of the real atmosphere. The foregoing suggests that the negative bias is due to the limitations of the CS3 surge model, which is common to both the HadGEM3-driven and the ERA-interim-driven results, and that the HadGEM3 simulation of storminess is comparable (at least by this metric) to the ERA reanalysis. In short, the atmospheric model is adequate, but we need a better surge model.

Further shape parameter results are shown in appendix D.

390 Very recently, Horsburgh et al. (2021) have simulated surge events which are higher than the CFB2018 best estimate of 1000-year return level. This finding is not contradictory to ours. Horsburgh et al. (2021) seek to identify unprecedented events which are *possible* in the present-day climate, without seeking to quantify their probability. Our focus is on improving the quantification of the probability of unprecedented events.

5 Sheerness case-study experiments

395 In view of the societal and economic importance of the Thames estuary, we further investigate the behaviour at Sheerness.

5.1 Tide-Surge timing

Figure 4 shows results of experiments making small (less than 24 hours) timing shifts in the phase relationship between the atmospheric forcing and the tide. We chose a spring tide and shifted the timing of the event relative to the tide. The curve labelled “SWL 0” shows the SWL of the shift which gives the maximum skew surge. The curve labelled “SWL 4” shows the
400 SWL when the event is shifted 4 hours later relative to the tide. The skew surge on the initial high tide (about nominal hour 24 in the figure) is reduced by about 1.2 metres. The overall maximum SWL (which occurs on the next high tide in the shifted case) is reduced by about 0.8 metres.

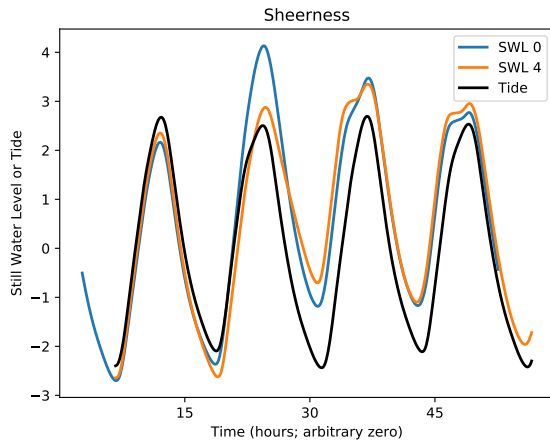


Figure 4. A simulated extreme event on a spring tide showing the effect of a shift of 4 hours in the timing of the event relative to the tide.

Clearly, a potentially extreme event may not be realised as an extreme SWL if does not happen to be in a conducive timing relationship with the tide. From a coastal defence viewpoint this is good, as it reduces the number of extreme SWLs which are realised. But from the viewpoint of identifying extreme events in a long model simulation it is a nuisance, because it can mean that potentially extreme events are hidden. To overcome this we performed a further simulation with the surge model in surge-only mode (see §6.1). In this mode no astronomical tides are included, and therefore all potentially extreme atmospheric events are realised as a surge.

5.2 Skew Surge/Tide dependence at Sheerness

Work by Williams et al. (2016) has shown that any dependence of skew surge on predicted high water cannot be readily quantified in the observational record, due to the dominance (in the record) of the variability of atmospheric storms. This conclusion has led to the exploitation of an assumed independence of skew surge and predicted high water as part of the effort to estimate present-day still water return levels — the so-called skew surge joint probability method which is used by CFB2018 (although they do note that such independence is not applicable everywhere).

This independence can be tested in model simulations, by repeating the same atmospheric storm in different astronomical tidal conditions – for example at spring and neap tide. Williams et al. (2016) perform four experiments of this kind (see their supplementary material) using reanalysed real-world storm data. We extended that work using 16 of the most extreme forcing events (in the sense that they create an extreme surge at Sheerness) from our HadGEM3-GC3-MM control simulation. Results are shown in appendix E. Here we give an example of a single event.

The largest skew surge event at Sheerness in the HadGEM3-GC3-MM simulation happened to arrive on a neap tide. Figure 5 shows that when this event was moved to a spring tide, the skew surge was significantly attenuated, from about 2 metres in the

case of the neap tide to about 1.63 metres in the case of the spring tide. Williams et al. (2016, their supplementary material S5) also found attenuation at Sheerness in model simulations of four events.

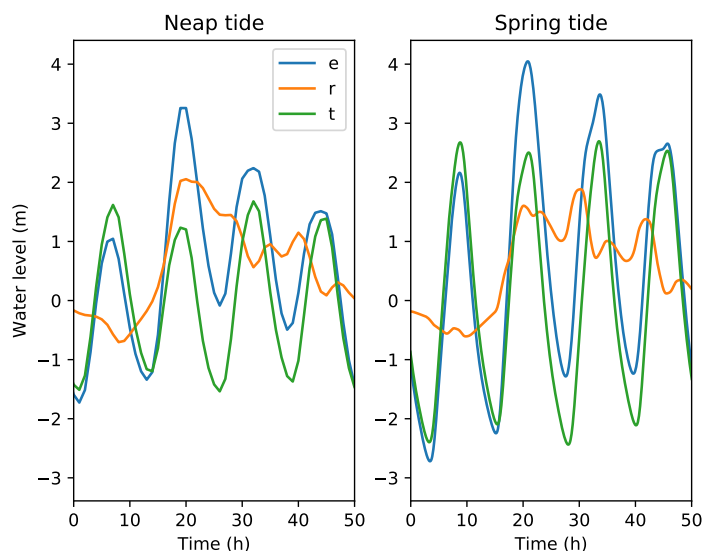


Figure 5. Left panel shows the largest skew surge in the HadGEM3-GC3-MM 483 year surge-and-tide simulation, which happened to arrive on a neap tide. Right panel shows the same atmospheric forcing applied to a spring tide. It can be seen that the realised skew-surge is dependent on the tide in this case. Key: e: still water elevation. t: astronomical tide. r: residual (i.e. e-t). X-axis shows time in hours with arbitrary zero.

Further skew-surge/tide dependence results are shown in appendix E. [D’Arcy et al. \(2021\) present observationally-based findings on tide-skew surge interaction.](#)

6 Sheerness: comparison of the most extreme simulated events with reconstructions of the 1953 event

6.1 Sheerness: Surge-only simulations

Our surge-only simulations are motivated by the sensitivity shown in Figs. 4 and 5, and discussed in §5. Using a numerical coastal shelf model it is possible to artificially eliminate the effect of the astronomical tide to create a surge-only simulation. Thus, issues of the timing relationship between surge and tide are eliminated in a surge-only simulation and so the sensitivity is avoided. This makes surge-only simulations well suited to comparing different sets of atmospheric forcing in terms of their surge-creation potential for a given location.

Figure 6 shows time series of water level at Sheerness for 16 events from our HadGEM3-GC3-MM surge-only simulation, in each case compared with a surge-only simulation driven by atmospheric data from a reconstruction (pers. comm: Erik van

435 Meijgaard, by email) of the 1953 storm using the KNMI/DMI limited area model RACMO (van Meijgaard et al., cited 2020). We selected first the largest 8 events in terms of the maximum value of the surge which they produced in surge-only mode. Compared to these events, the 1953 surge-only simulation has a conspicuous long duration. So we also sought events of long duration by convolving the surge only time series with the kernel shown in Fig. 6 panel i. Then we identified the 8 largest maxima in the convolved signal. All 16 events are independent (all separated from each other by at least a year). The kernel was designed to represent the important features of the RACMO-driven surge-only simulation, i.e. the approximate duration and shape of the time series plot. The purpose of convolution with the kernel is to identify those events which correlate well (in terms of their time series plot) with the RACMO-driven simulation, in other words, events which not only produce a large surge, but are also of comparable duration to the RACMO-driven simulation. The kernel was not used to modify events, but simply to identify significant ones.

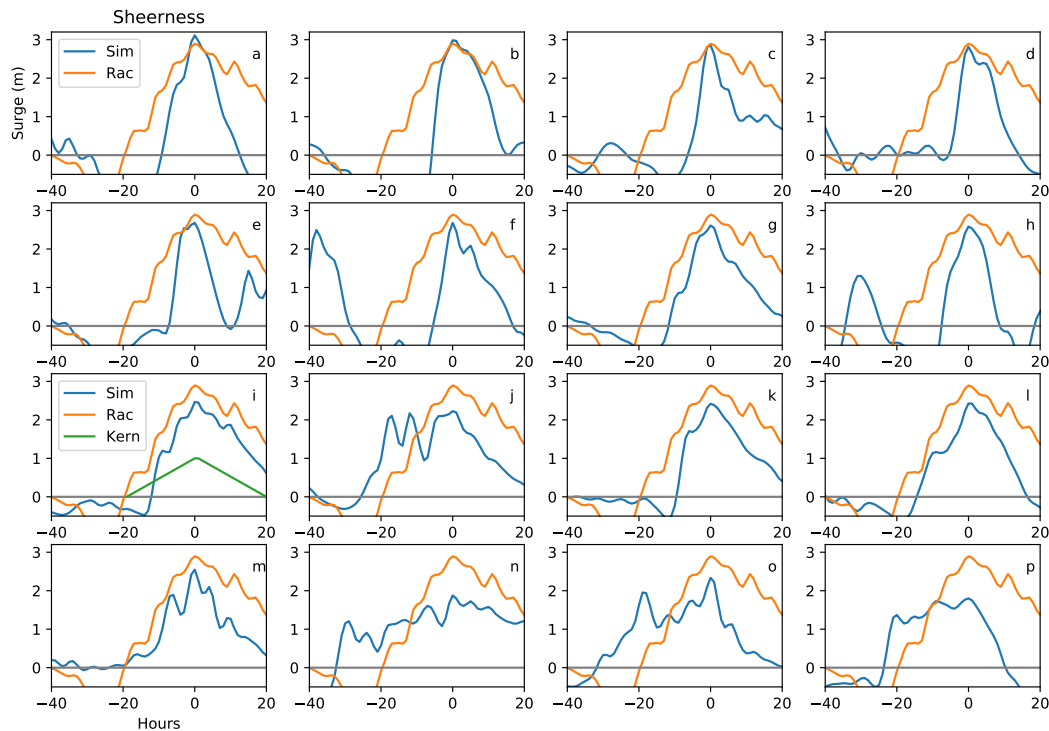


Figure 6. 16 events from the HadGEM3-GC3-MM surge-only simulation (“Sim”), in each case compared with the RACMO-driven surge only simulation (“Rac”). Panels a to h show the 8 largest independent surge-only events. Panels i to p show events which are both large and have substantial duration (the original time series before convolution is shown). X-axis is time in hours with arbitrary zero.

445 This shows that in the 483-year surge-only simulation

1. No simulated event exceeded the 1953 reconstruction in terms of both maximum surge *and* duration.
2. Two simulated events exceeded the 1953 reconstruction in terms of maximum surge, and several more were comparable.
3. Several simulated events were of comparable duration to the 1953 reconstruction, but exhibited a smaller maximum surge.

450 **6.2 Sheerness: Surge and Tide simulations**

Having used the surge-only mode to identify 16 potentially-extreme events in the HadGEM3-GC3-MM simulation, for each event we experimented with adjusting the timing of the event in a surge-and-tide simulation to maximise the skew surge realised. We did this twice: once for a spring tide and once for a neap tide. Figure 7 shows (bar “S”) the overall (i.e. over all 16 events) maximum skew surge realised on a spring tide and similarly the overall maximum skew surge realised on a neap tide (bar “N”). Figure 7 also shows (bar “H”) the maximum skew surge realised in the original HadGEM3-GC3-MM surge-and-tide simulation, in which the timings were not artificially adjusted, so that the surge/tide phase relationship was essentially random (as in the real world). For reference an extreme (entirely artificial) case is shown (bar “Z”) in which no tidal forcing is included (see §6.1). In reality, of course, the tide is always present. Wadey et al. (2015) tabulate estimates of high water level at Sheerness for the 1953 event from four different sources. They also give a best estimate of 4.74 metres, and a corresponding best-estimate skew surge of 2.16 metres. This implies an astronomical tide of $4.74 - 2.16 = 2.58$ metres. Thus, to obtain the four skew surge estimates labelled “W” in Fig. 7, we subtract this tide from each of their four tabulated estimates of high water level.

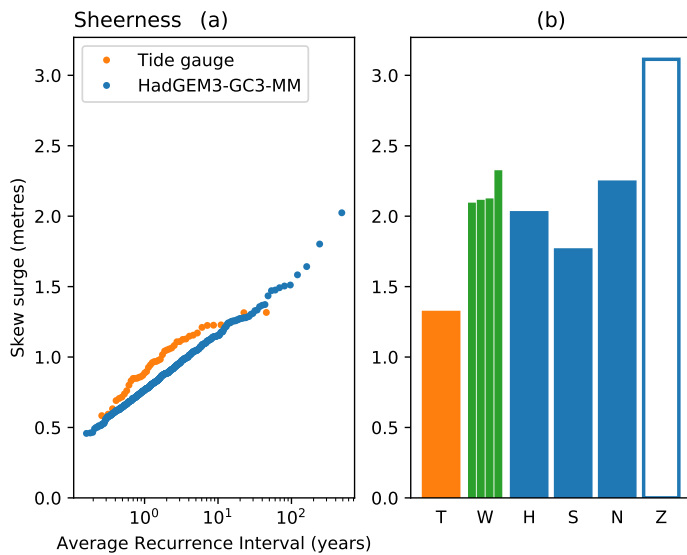


Figure 7. Observed/estimated (in orange) and modelled (in blue) skew surges at Sheerness. (a): Empirical return level plot showing annual maxima of skew surge from tide gauge data as used in CFB2018 (orange) and from 483-year model simulation (blue). Years with less than 75% of available data are excluded from the tide gauge data analysis. (b): Skew surge maxima. Key: T(ide): Tide gauge (max from panel (a)). W(adey): Data from four different sources for the 1953 event, as tabulated by Wadey et al. (2015). H(adGEM3): Model max skew. Phase relationship between atmospheric events and astronomical tide is essentially random over the 483 year simulation. S(pring): Model max skew when events are artificially shifted to coincide with Spring tide. N(eap): Model max skew when events are artificially shifted to coincide with Neap tide. Z(ero): Model max skew when astronomical tides are excluded (surge-only simulation: “Zero tide”).

Figure 7 shows that the strongest atmospheric forcing in the model simulation can produce a skew surge which is comparable to estimates of the 1953 skew surge at Sheerness. The largest skew surge (bar “H”) in the simulation in which the timings of atmospheric forcing and tides are essentially random, lies just below the range of skew surge estimates based on data tabulated by Wadey et al. (2015). The largest spring-tide skew surge (bar “S”, when the timings of atmospheric forcing are adjusted so that the atmospheric events coincide with a spring tide) is smaller than the observational estimates, due to the surge-tide interaction at this site. The largest neap-tide skew surge (bar “N”, when the timings of atmospheric forcing are adjusted so that the atmospheric events coincide with a neap tide) lies within the range of skew surge estimates based on data tabulated by Wadey et al. (2015).

~~We do not attempt to explicitly quantify a return period of the 1953 SWL at Sheerness because an evaluation of the probability distribution of SWL requires a convolution of the distributions of skew surge and of tide (see §3.3). CS3, in common with other shelf models, is known to exhibit tidal errors, typically under-predicting the range. However, the fact that the 483-year surge-only simulation produces more than one event of comparable magnitude to the simulated 1953 event~~

475 ~~suggests that the return period of the 1953 atmospheric forcing is less than 483 years, in so far as the model is realistic.~~
~~Wadey et al. (2015) suggest a return period of 429 years for the surge event at Sheerness.~~

7 Summary and Conclusions

HadGEM3-GC3-MM is a state-of-the-art global climate model of the CMIP6 generation. Modifications including the ENDGame revision to the dynamical core have been shown to increase synoptic variability (Williams et al., 2015), improving the representation of the storm tracks compared to HadGEM2-AO (the Hadley Centre model which contributed to CMIP5). We have
480 shown that a 483-year control simulation of HadGEM3-GC3-MM, in combination with a barotropic storm surge model of the north west European coastal shelf, is capable of directly simulating realistic extreme storm surges for some sites around the UK coastline, as evaluated against observations (Fig. 1). In particular, our modelling system simulates several surge events at Sheerness (on the Thames Estuary) which are comparable to best estimates of the catastrophic 1953 storm (Figs. 6 and 7).

485 We extend the skew surge–tide dependence results of Williams et al. (2016). Our simulations suggest that skew surge–tide dependence can have a substantial effect on the most extreme surges at Sheerness (Fig. 7 and appendix E).

Furthermore, around the whole of the UK coastline we find that the spatial pattern of variations in the three parameters which describe the extreme tail of the storm surge distribution is very well reproduced by the simulation (Fig. 2). In particular, the observed spatial variations in the shape parameter are reproduced by the simulation. This is important because

- 490
- it gives further credibility to both diagnoses of the spatial variations
 - the shape parameter is the main source of uncertainty in estimates of unprecedented events (appendix C)
 - the length of the simulation (much greater than the length of the observational record) helps to constrain the shape parameter ~~with less subjectivity (§4.1).~~

A typical simulated shape parameter for an individual site is more negative than (but within the uncertainty of) that diagnosed
495 by CFB2018 (Fig. 2 panel (d)). This negativity arises at a wide spread of sites. Such spatial uniformity of the negativity strongly suggests an underlying difference rather than a chance/sampling difference. Sub-sampling the simulation with sample sizes matching the tide gauge record lengths supports that suggestion and shows that our model shape parameters are biased low relative to those diagnosed from tide-gauge observations. However, that is also the case when our surge model is driven by a good quality atmospheric reanalysis, suggesting that the bias comes from shortcomings at the surge modelling stage rather
500 than the atmospheric forcing.

We conclude, then, that our atmospheric model, HadGEM3-GC3-MM, has the potential to help constrain estimates of unprecedented UK storm surges, but that improvements at the surge modelling stage are required.

8 Suggestions for Further Work

505 The model/observation departures seen in Fig. B1 have a systematic feature which is consistent over spatial regions, e.g., south-west and north-west UK. This suggests that it should be possible to account for these departures through a smooth spatial function which maps the differences in quantiles between the observations and model data. With this adjustment it is possible that the currently identified under-estimation may be corrected before making the tail-based GEV/GPD comparisons shown here.

510 The shape parameter estimates in Fig. 2(d) show some site-to-site variations which are more pronounced than the broader smooth variations across coastlines. This suggests that they would also benefit from the penalty-based approach used by CFB2018.

One advantage of modelled data is the ability to give estimates at ungauged sites. We have not exploited this ability here, but we anticipate that it will form the basis of further work.

515 *Data availability.* The tide gauge data used in the CFB2018 report are available to download from the National Tidal and Sea Level Facility (ntslf.org). The CFB2018 shape parameters can be read from their figure E.1 (Environment Agency, 2018). The CFB2018 GEV scale parameters as shown in Fig. 2(b), in metres, (for sites Newlyn... Portrush in the order shown on the X-axis of Fig. 2) are:

0.0835, 0.0733, 0.0847, 0.1031, 0.1369, 0.1790, 0.1574, 0.1484, 0.1140, 0.0940, 0.1639, 0.1063, 0.1161, 0.1305, 0.1228, 0.1815, 0.1593, 0.1358, 0.1757, 0.1490, 0.1471, 0.1202, 0.0955, 0.1368, 0.0684, 0.0913, 0.0922, 0.0966, 0.1049, 0.1178, 0.1333, 0.1618, 0.1748, 0.2153, 0.1715, 0.1629, 0.1557, 0.1091, 0.0985, 0.0937, 0.0948, 0.1201, 0.0938, 0.1273

520 Simulated sea levels at the tide gauge sites as used in our analysis are available from the first author on request. All of the analysis was undertaken using the open source languages R and Python.

Appendix A: Surge model grid and tide gauge locations

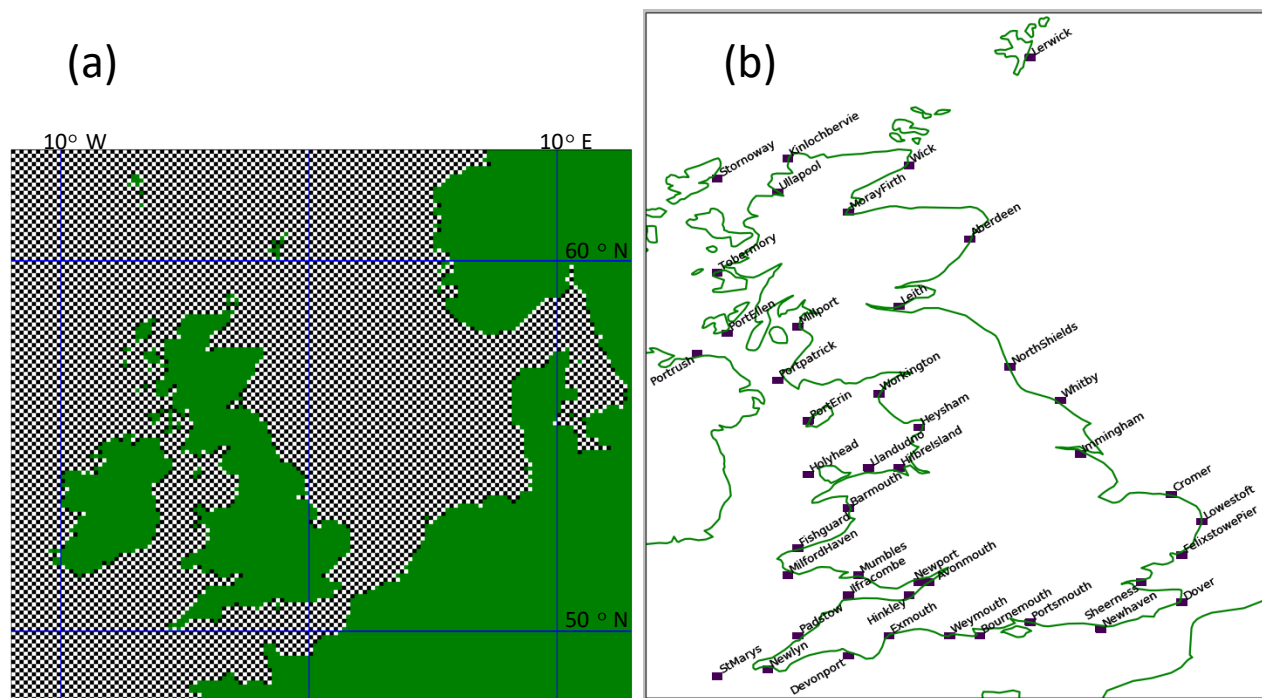


Figure A1. (a) Domain and grid of the CS3 coastal shelf model. Grid size is 1/9 degree in latitude and 1/6 degree in longitude, which results in near-square grid cells at the latitude of the UK. (b) Tide gauge locations.

Appendix B: Empirical return level plots for UK tide gauges

Figure B1 shows empirical return level plots for 44 tide gauge locations around the UK. For simplicity we use annual maxima only. The observational annual maxima are limited to years in which the tide gauge data is at least 75 % complete. Plotting positions are evaluated using the Weibull formula (Weibull, 1939). Working from left to right along the rows (and then downwards through the rows as in reading), the sequence of plots follows the clockwise sequence of Fig. 2 as described in §4.1.

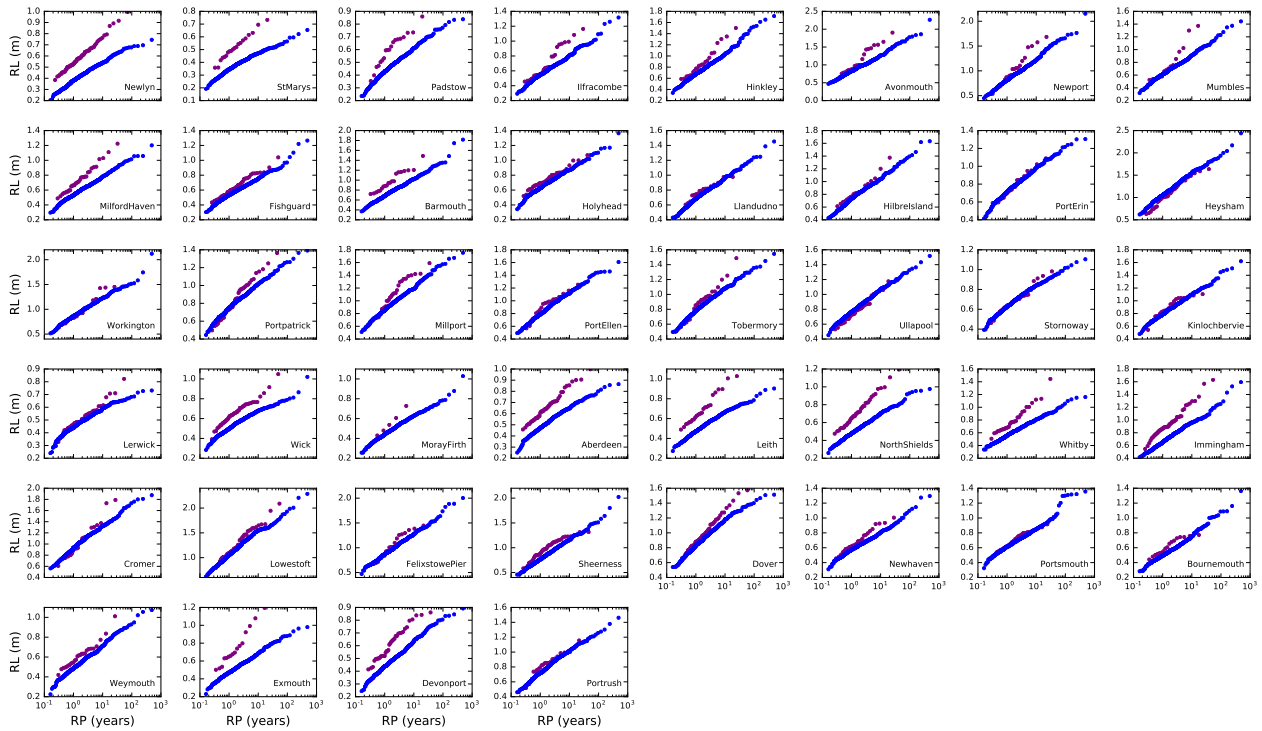


Figure B1. Empirical return level plots for 44 UK tide gauges. Blue-Purple shows observational (tide-gauge) data. Purple-Blue shows data from the 483-year HadGEM3-GC3-MM simulation.

530 It-[Figure B1](#) shows some major departures between the model and observed data across the distribution of skew surges, but particularly in the tails. The model does not give higher quantiles than the observed data at any site. However, it can be seen that at some locations the model produces a plausible simulation of the observed return level plot and a plausible extrapolation of the return level plot to return periods outside of the observational record. This is discussed further in §4.

Appendix C: Shape-parameter uncertainty dominance

For short record lengths, unconstrained maximum-likelihood estimation is known to give “noisy” and implausible shape parameters (Coles and Dixon, 1999; Martins and Stedinger, 2000), see also §3.3. We illustrate this in Fig. C1 with a GEV fit to
 535 tide-gauge data.

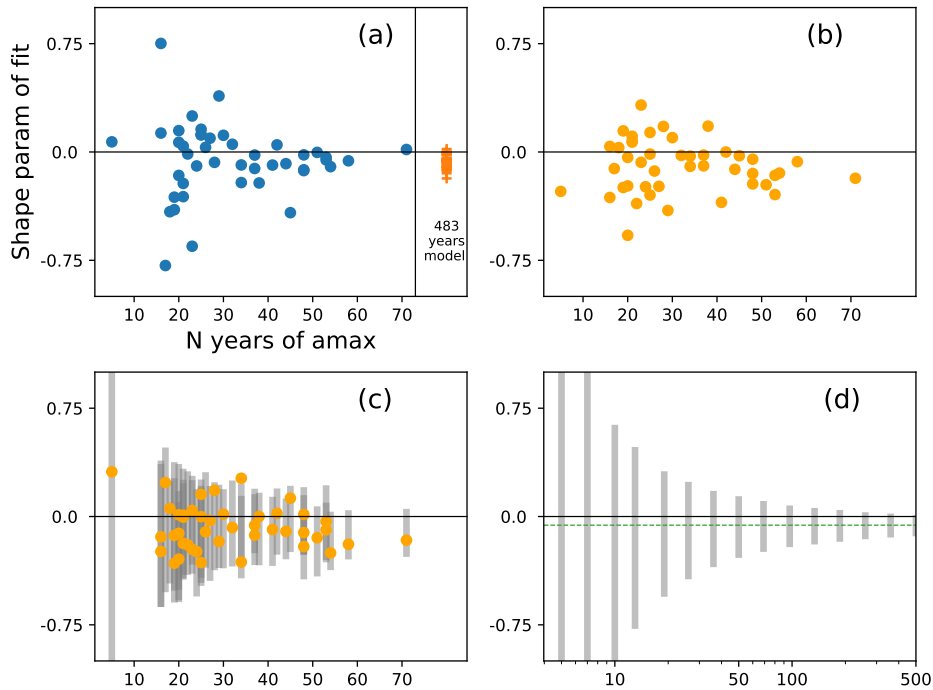


Figure C1. Short record lengths lead to noisy MLE shape-parameter estimates. (a) Shape parameter of GEV fit to tide gauge data against the number of annual maxima fitted (blue dots, one for each tide gauge), and shape parameter of GEV fit to model data (orange crosses, one for each tide gauge, all 483 years). Note that the range of fitted shape parameters reduces (“tapers”) as record length increases. (b) Model data for each port is cut down to a (random) sub-sample having the same length as the observational record at that port and then fitted in the same way as the observations. A similar tapering pattern emerges. (c) as (b) but a different random sub-sample. (d) as (b) but a different random sub-sample. Also shown (grey) is the 5 to 95 percentile range from one hundred such sub-samples at each port. (d) Similar to (a, b, c) but using pseudo-random variates drawn from a GEV distribution with scale and shape parameter which are typical of the values found around the UK (0.12 metres and -0.06 respectively). Sample size varies from 5 to 500, as shown by the logarithmically-scaled x-axis. For each sample, GEV parameters are fitted by maximum likelihood estimation. For each sample size, 2000 samples are drawn and the 5 to 95 percentile of the fitted shape parameters is shown by the grey bar. The dashed green line shows a shape of -0.06. The 5 to 95 percentile intervals for sample sizes 5, 7, and 10 exceed the Y-axis limits.

In similar plots for the location and for the scale parameter, no such tapering is exhibited. In this illustration, tide gauge records of length greater than about 40 years have a fitted shape parameter which is within the range of the model fitted shape parameters; only in short records are large positive or negative shape parameters found. Figure C1 (b), (c) and (d) confirm that the tapering is a result of record length.

540 Associated with this, for observational record lengths, the uncertainty in the shape parameter dominates the uncertainty in inferred return levels for long return periods. This is illustrated in Fig. C2.

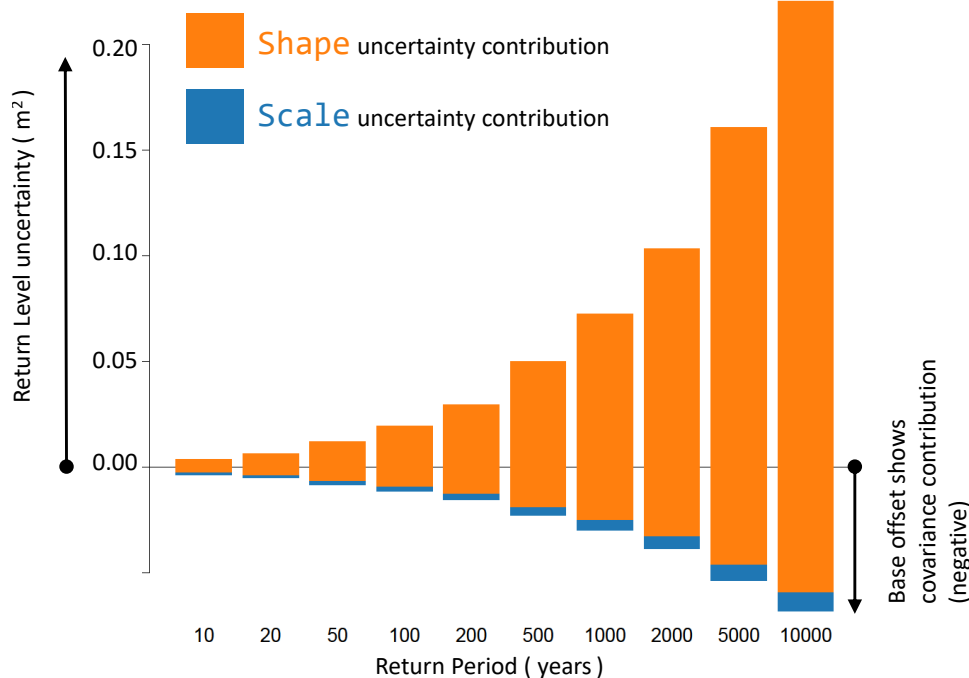


Figure C2. The uncertainty (variance) in different return levels is partitioned into contributions from uncertainty in the shape parameter, uncertainty in the scale parameter, and a negative contribution from the covariance of these two parameters, shown as an offset to the base of the bars. Uncertainty in the shape parameter becomes dominant at long return periods.

Figure C2 shows the sources of uncertainty in return level for ten different return periods. The data were constructed as follows. We took representative shape and scale parameters (we used the CFB2018 parameters for skew surge at Sheerness) and a representative record length of 45 years. We simulated 45 years of 94 % threshold exceedance data by inverse transform sampling. We estimated the (GPD) parameters of the sample by maximum likelihood estimation without any prior constraint. For each of ten return levels, we estimated the uncertainty using the delta method (e.g. Coles, 2001) as follows.

$$\text{Var}(R) = \frac{\partial R}{\partial \tilde{\sigma}} \text{Var}(\tilde{\sigma}) \frac{\partial R}{\partial \tilde{\sigma}} + \frac{\partial R}{\partial \xi} \text{Var}(\xi) \frac{\partial R}{\partial \xi} + 2 \frac{\partial R}{\partial \tilde{\sigma}} \text{Cov}(\tilde{\sigma}, \xi) \frac{\partial R}{\partial \xi} \quad (\text{C1})$$

where R is return level, $\tilde{\sigma}$ is the GPD scale parameter, and ξ is the shape parameter. The variance and covariance terms, which are determined from the curvature of the likelihood surface, are evaluated during the likelihood maximisation routine. The three terms on the right-hand side of equation C1 are the contributions to the return level uncertainty from the GPD scale parameter uncertainty, the shape parameter uncertainty, and the covariance of the two parameters, respectively. To increase confidence in the uncertainty estimates we repeated the sampling many times and averaged over each contribution. A further contribution to uncertainty is the choice of threshold, but this contribution is usually found to be small (Coles, 2001) and is neglected here. We tested some alternative approaches (not shown here), for example fitting a GEVD to the annual maxima instead of GPD to POT. The essential result — the dominance of the shape parameter uncertainty — is robust and was not affected by

the use of alternative approaches. The result holds for all of the nine locations that we tried: Newlyn, Fishguard, Holyhead, Stornoway, Lerwick, Aberdeen, Cromer, Lowestoft and Sheerness. In each case we simulated a record length corresponding to the available tide gauge data for that location.

Appendix D: Further shape parameter results

560 Figure 3 in the main text shows results of fitting the GPD distribution to peaks over a threshold. Here, in Fig. D1, we extend that figure to include results of fitting the GEV distribution to annual maxima.

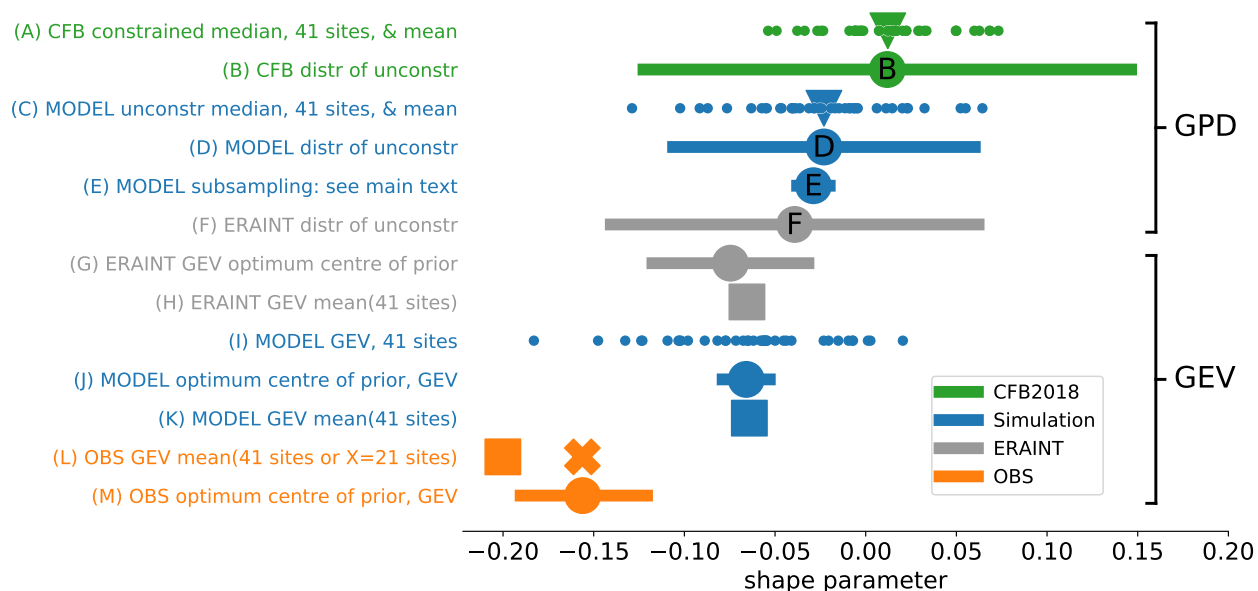


Figure D1. Further shape parameter results. (A to F): diagnosed by GPD fit to POT as Fig. 3. (G to M): diagnosed by GEV fit to annual maxima.

To probe the model–tide gauge shape parameter bias further we performed some more experiments. We made an unconstrained GEV fit to tide-gauge annual maxima of skew surge at each of the 41 sites. Owing to the short record length the results are very noisy and include some implausible values. Nevertheless the mean of these 41 results is shown by the filled square on line (L). The mean for the 21 sites with the longest record lengths is also shown, by the cross on line (L). For comparison, the mean of the GEV fit to simulated data (41 sites) is shown in line (K). The GEV fit to simulated data for each of the 41 sites is shown in line (I) (i.e. the square on line (K) is the mean of the points on line (I)). To compensate for the short observational record length we experimented with applying a prior to the shape parameter of the GEV fit to the observed annual maxima of skew surge. Following CFB2018, we used a normal prior with a standard deviation of 0.0343, but we varied the centre (i.e. the

565

570 mean) of the prior. For each site, we produced a maximum likelihood fit by maximising the log-likelihood in the usual way. We summed log-likelihoods over all sites to give an overall log-likelihood for that value of the centre of the prior. The value of the centre of the prior which maximised the log-likelihood is shown by the disc in line (M). Standard techniques (Coles, 2001) enable us to identify a 95 % confidence interval, shown by the solid straight line in line (M). We applied the same approach to the simulated skew surges to give the results shown in line (J).

575 We did some further shape parameter evaluations with surges generated by CS3 driven by the ERA-interim atmospheric reanalysis (Dee et al., 2011) that has been downscaled with the Swedish Meteorological and Hydrological Institute (SMHI) Rossby Centre regional atmospheric model (RCA4) as part of the Euro-CORDEX experiment (Jacob et al., 2014). The ERA-interim GEV-mean (corresponding to the square in line (L)) is shown in line (H). The GEV-optimum centre of prior (corresponding to line (M)) is shown in line (G). In all cases, our ERA-interim-based shape parameter results are in better agreement
580 with our model results than with the CFB shape parameter results. As discussed in the main text, this suggests that the bias is not a result of any deficiency in the climate-model compared to the ERA-interim reanalysis, and that it is more likely arising from the continental shelf modelling stage.

The GEV fits to skew surge data give more negative shape parameter results than GPD fits. This is surprising since both methods (GEV and GPD) are asymptotically unbiased (Coles, 2001). We further experimented with fitting to some entirely
585 artificial data (generated pseudo-random numbers) of comparable size to the simulations. We did not find any evidence of a consistent negative shape bias in the MLE when using a GEV compared to a GPD fit for any distribution of pseudo-random numbers that we tested (including Gumbel, normal, and GEV with +ve and -ve shape parameter). Thus the GEV vs GPD bias may be hinting at a departure from the conditions which are required for accurate statistical modelling of the extremes.

Appendix E: Further skew-surge/tide dependence results at Sheerness

590 For each of the 16 extreme events shown in Fig. 6 we re-ran the simulation, adjusting the timing such that the event coincided with a spring and then a neap tide. In each case we also made small timing adjustments to realise the maximum skew surge. In all cases the skew surge was attenuated on the spring tide relative to the neap tide. The size of the skew surge, and the spring-neap difference in the skew surge are shown in Fig. E1.

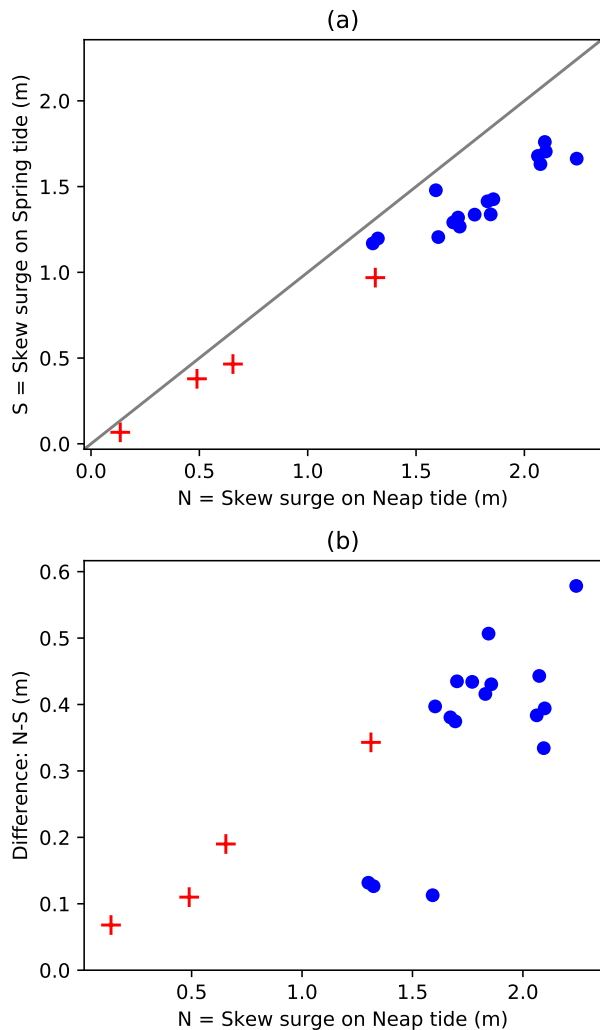


Figure E1. (a) Skew surge realised on spring tide (“S”, Y-axis) vs. skew surge realised on neap tide (“N”, X-axis) (b) Difference (N-S) vs. N. Blue dots show our experiments using atmospheric forcing from the 16 most extreme Sheerness events in the HadGEM3-GC3-MM simulation. Red crosses show data from Williams et al.(2016).

As discussed in the main text, such experiments have been conducted before by Williams et al. (2016); their results are shown
 595 by the red crosses in Fig. E1. Williams et al. (2016) also show scatter plots of observed skew surge and tide. In Fig. E2 we
 show our model results overlain on a reproduction of the Sheerness panel from their figure S2. The blue points show simulated
 skew surge against simulated astronomical high water for the 16 extreme atmospheric events with timing adjusted such that the
 event coincided with a simulated spring tide and a simulated neap tide, as described above. The artificial case of no tide is also
 shown. Grey points are as in Williams et al. (2016). It can be seen that our model results do not look out of context compared to

600 the observations in terms of the negative correlation. For reference the extreme (entirely artificial) case of simulated tide-surge interaction is also shown, in which no tidal forcing is included (see §6.1). In reality, of course, the tide is always present.

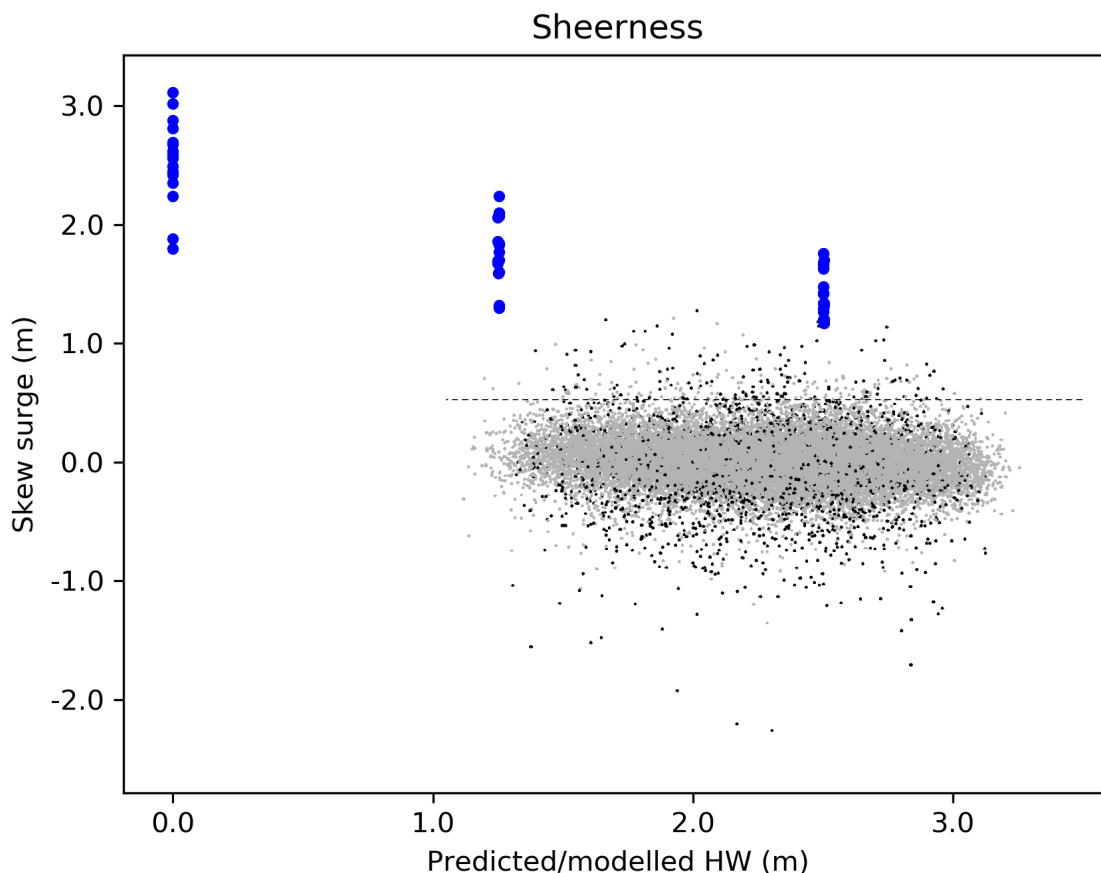


Figure E2. Our skew-surge/tide interaction results for Sheerness overlain on a reproduction of Williams et al.(2016), their figure S2. The results for a modelled high water of zero are included for comparison. They come from a very artificial simulation in which no tidal forcing is included. In reality, of course, the tide is always present.

Author contributions. TH devised the experiment, performed most of the new analysis, and wrote the article. SDW performed the original CFB2018 analysis on the tide gauge data and replicated it for the simulation.

Competing interests. The authors declare no competing interests.

605 *Acknowledgements.* This work was supported by the Met Office Hadley Centre Climate Programme funded by BEIS and Defra. [We thank Andreas Sterl, Eleanor D'Arcy, and Jon Tawn for their generous reviews, which helped to improve the initial submitted manuscript of this work.](#) Special thanks to Jenny Sansom for helpful discussions and for providing some of the data used in the evaluation. Thanks to Erik van Meijgaard ~~and Andreas Sterl~~ at KNMI for comments on an early draft and for permission to use their RACMO-based atmospheric reconstruction of the 1953 event, and to Mark Pickering at National Oceanography Centre Southampton for helping to supply the data.

610 Thanks to Graham Siggers (HR Wallingford) and Matt Palmer (Met Office Hadley Centre) for helpful comments on an early draft of this paper. Thanks to Jeff Ridley for help with the HadGEM3-GC3-MM data. Thank you to ~~Jon Tawn~~, Simon Brown and Rob Shooter for [helpful discussions and](#) guidance with extreme value modelling. This publication contains public sector information licensed under the Open Government Licence v3.0.

References

- 615 Arns, A., Wahl, T., Haigh, I., Jensen, J., and Pattiaratchi, C.: Estimating extreme water level probabilities: A comparison of the direct methods and recommendations for best practise, *Coastal Engineering*, 81, 51–66, 2013.
- Batstone, C., Lawless, M., Tawn, J., Horsburgh, K., Blackman, D., McMillan, A., Worth, D., Laeger, S., and Hunt, T.: A UK best-practice approach for extreme sea-level analysis along complex topographic coastlines, *Ocean Engineering*, 71, 28–39, 2013.
- Bernier, N. and Thompson, K.: Predicting the frequency of storm surges and extreme sea levels in the northwest Atlantic, *Journal of Geophysical Research: Oceans*, 111, 2006.
- 620 Brown, J. M., Souza, A. J., and Wolf, J.: Surge modelling in the eastern Irish Sea: present and future storm impact, *Ocean Dynamics*, 60, 227–236, 2010.
- Coles, S.: An introduction to statistical modeling of extreme values, Springer, 2001.
- Coles, S. G. and Dixon, M. J.: Likelihood-based inference for extreme value models, *Extremes*, 2, 5–23, 1999.
- 625 D’Arcy, E., Tawn, J., Joly-Laugel, A., and Sifnioti, D. E.: Accounting for seasonality in extreme sea level estimation, (in preparation), 2021.
- de Vries, H., Breton, M., de Mulder, T., Krestenitis, Y., Ozer, J., Proctor, R., Ruddick, K., Salomon, J. C., and Voorrips, A.: A comparison of 2D storm surge models applied to three shallow European seas, *Environmental Software*, 10, 23 – 42, 1995.
- de Winter, R. C., Sterl, A., and Ruessink, B. G.: Wind extremes in the North Sea Basin under climate change: An ensemble study of 12 CMIP5 GCMs, *J. Geophys. Res. Atmos.*, 118, 1601–1612, <https://doi.org/10.1002/jgrd.50147>, 2013.
- 630 Dee, D. P., Uppala, S., Simmons, A., Berrisford, P., Poli, P., Kobayashi, S., Andrae, U., Balmaseda, M., Balsamo, G., Bauer, d. P., et al.: The ERA-Interim reanalysis: Configuration and performance of the data assimilation system, *Quarterly Journal of the Royal Meteorological Society*, 137, 553–597, 2011.
- Edwards, T.: Future of the sea: impacts of sea level rise on the UK, Government Office for Science, Foresight Future of the sea, Government Office for Science., 2017.
- 635 Environment Agency: Coastal Flood Boundary Conditions for the UK: update 2018, Tech. rep., <https://environment.data.gov.uk/dataset/84a5c7c0-d465-11e4-b0bd-f0def148f590>, Available to download from <https://environment.data.gov.uk>, 2018.
- Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., and Taylor, K. E.: Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6) experimental design and organization, *Geoscientific Model Development*, 9, 1937–1958, 2016.
- Flather, R. A.: A Storm Surge Prediction Model for the Northern Bay of Bengal with Application to the Cyclone Disaster in April 1991, *J. Phys. Oceanogr.*, 24, 172–190, 1994.
- 640 Flather, R. A.: Existing operational oceanography, *Coastal Engineering*, 41, 13–40, 2000.
- Furner, R., Williams, J., Horsburgh, K., , and Saulter, A.: NEMO-surge: Setting up an accurate tidal model., Forecasting Research Technical Report No. FRTR610, Met Office, UK., https://digital.nmla.metoffice.gov.uk/IO_0499462b-eea6-4313-8761-ee41d533872c/, 2016.
- Grabemann, I., Gaslikova, L., Brodhagen, T., and Rudolph, E.: Extreme storm tides in the German Bight (North Sea) and their potential for amplification, *Natural Hazards and Earth System Sciences*, 20, 1985–2000, <https://doi.org/10.5194/nhess-20-1985-2020>, 2020.
- 645 Haigh, I. D., Nicholls, R., and Wells, N.: A comparison of the main methods for estimating probabilities of extreme still water levels, *Coastal Engineering*, 57, 838–849, 2010.
- Haigh, I. D., Wadey, M. P., Wahl, T., Ozsoy, O., Nicholls, R. J., Brown, J. M., Horsburgh, K., and Gouldby, B.: Spatial and temporal analysis of extreme sea level and storm surge events around the coastline of the UK, *Scientific data*, 3, 1–14, 2016.

- 650 Haigh, I. D., Ozsoy, O., Wadey, M. P., Nicholls, R. J., Gallop, S. L., Wahl, T., and Brown, J. M.: An improved database of coastal flooding in the United Kingdom from 1915 to 2016, *Scientific data*, 4, 170 100, 2017.
- Horsburgh, K., Williams, J., Flowerdew, J., and Mylne, K.: Aspects of operational forecast model skill during an extreme storm surge event, *Journal of Flood Risk Management*, 1, 213–221, 2008.
- Horsburgh, K., Haigh, I., Williams, J., De Dominicis, M., Wolf, J., Inayatillah, A., and Byrne, D.: “Grey swan” storm surges pose a greater
655 coastal flood hazard than climate change, accepted for publication in *Ocean Dynamics*, 2021.
- Howard, T. and Clark, P.: Correction and downscaling of NWP wind speed forecasts, *Meteorological Applications: A journal of forecasting, practical applications, training techniques and modelling*, 14, 105–116, 2007.
- Howard, T. and Palmer, M. D.: Sea-level rise allowances for the UK, *Environmental Research Communications*, 2, 035 003, <https://doi.org/10.1088/2515-7620/ab7cb4>, 2020.
- 660 Howard, T., Palmer, M. D., and Bricheno, L. M.: Contributions to 21st century projections of extreme sea-level change around the UK, *Environmental Research Communications*, 1, 095 002, <https://doi.org/10.1088/2515-7620/ab42d7>, 2019.
- Jacob, D., Petersen, J., Eggert, B., Alias, A., Christensen, O. B., Bouwer, L. M., Braun, A., Colette, A., Déqué, M., Georgievski, G., Georgopoulou, E., Gobiet, A., Menut, L., Nikulin, G., Haensler, A., Hempelmann, N., Jones, C., Keuler, K., Kovats, S., Kröner, N., Kotlarski, S., Kriegsmann, A., Martin, E., van Meijgaard, E., Moseley, C., Pfeifer, S., Preuschmann, S., Radermacher, C., Radtke, K.,
665 Rechid, D., Rounsevell, M., Samuelsson, P., Somot, S., Soussana, J. F., Teichmann, C., Valentini, R., Vautard, R., Weber, B., and Yiou, P.: EURO-CORDEX: New high-resolution climate change projections for European impact research, *Regional Environmental Change*, <https://doi.org/10.1007/s10113-013-0499-2>, 2014.
- Lowe, J., Howard, T., Pardaens, A., Tinker, J., Holt, J., Wakelin, S., Milne, G., Leake, J., Wolf, J., Horsburgh, K., Reeder, T., Jenkins, G., Ridley, J., Dye, S., and Bradley, S.: UK Climate Projections Science Report: Marine and Coastal Projections, Met Office Hadley Centre,
670 Exeter, UK, 2009.
- Martins, E. S. and Stedinger, J. R.: Generalized maximum-likelihood generalized extreme-value quantile estimators for hydrologic data, *Water Resources Research*, 36, 737–744, 2000.
- Menéndez, M. and Woodworth, P. L.: Changes in extreme high water levels based on a quasi-global tide-gauge data set, *Journal of Geophysical Research*, 115, <http://dx.doi.org/10.1029/2009jc005997>, 2010.
- 675 O’Neill, C., Saulter, A., Williams, J., and Horsburgh, K.: NEMO-surge: Application of atmospheric forcing and surge evaluation, *Forecasting Research Technical Report No. FRTR619*, Met Office, UK., https://digital.nmla.metoffice.gov.uk/IO_53fa4f69-432c-40bb-9481-8c7dfbd6492d/, 2016.
- Palmer, M., Howard, T., Tinker, J., Lowe, J., Bricheno, L., Calvert, D., Edwards, T., Gregory, J., Harris, G., Krijnen, J., and Roberts, C.: UK Climate Projections Science Report: Marine and Coastal Projections, Met Office Hadley Centre, Exeter, UK, 2018.
- 680 Palmer, M., Gregory, J. M., Bagge, M., Calvert, D., Hagedoorn, J., Howard, T., Klemann, V., Lowe, J., Roberts, C., Slangen, A., et al.: Exploring the drivers of global and local sea-level change over the 21st century and beyond, *Earth’s Future*, 8, e2019EF001 413, 2020.
- Pörtner, H.-O., Roberts, D. C., Masson-Delmotte, V., Zhai, P., Tignor, M., Poloczanska, E., Mintenbeck, K., Nicolai, M., Okem, A., Petzold, J., et al.: IPCC Special Report on the Ocean and Cryosphere in a Changing Climate, IPCC Intergovernmental Panel on Climate Change: Geneva, Switzerland, 1, 2019.
- 685 Priestley, M. D., Ackerley, D., Catto, J. L., Hodges, K. I., McDonald, R. E., and Lee, R. W.: An overview of the extratropical storm tracks in CMIP6 historical simulations, *Journal of Climate*, 33, 6315–6343, 2020.

- Roberts, M., Vidale, P., Senior, C., Hewitt, H., Bates, C., Berthou, S., Chang, P., Christensen, H., Danilov, S., Demory, M.-E., et al.: The benefits of global high resolution for climate simulation: process understanding and the enabling of stakeholder decisions at the regional scale, *Bulletin of the American Meteorological Society*, 99, 2341–2359, 2018.
- 690 Shaw, T., Baldwin, M., Barnes, E. A., Caballero, R., Garfinkel, C., Hwang, Y.-T., Li, C., O’Gorman, P., Rivière, G., Simpson, I., et al.: Storm track processes and the opposing influences of climate change, *Nature Geoscience*, 9, 656–664, 2016.
- Shepherd, T. G.: Atmospheric circulation as a source of uncertainty in climate change projections, *Nature Geoscience*, 7, 703–708, 2014.
- Sterl, A., Van den Brink, H., de Vries, H., Haarsma, R., and van Meijgaard, E.: An ensemble study of extreme storm surge related water levels in the North Sea in a changing climate, *Ocean Science*, 5, 369–378, <https://doi.org/10.5194/os-5-369-2009>, 2009.
- 695 Tawn, J.: Estimating probabilities of extreme sea-levels, *Appl. Statist.*, 41, 77–93, 1992.
- Taylor, K., Stouffer, R., and Meehl, G.: An Overview of CMIP5 and the experiment design., *Bull. Amer. Meteor. Soc.*, 93, 485–498, <https://doi.org/10.1175/BAMS-D-11-00094.1>, 2012.
- Thompson, V., Dunstone, N. J., Scaife, A. A., Smith, D. M., Slingo, J. M., Brown, S., and Belcher, S. E.: High risk of unprecedented UK rainfall in the current climate, *Nature communications*, 8, 1–6, 2017.
- 700 Van den Brink, H., Können, G., Opsteegh, J., van Oldenborgh, G. J., and Burgers, G.: Improving 10⁴-year surge level estimates using data of the ECMWF seasonal prediction system, *Geophysical Research Letters*, 31, 2004.
- van Meijgaard, E., van Uft, L., van de Berg, W., Bosveld, F. C., van den Hurk, B., Lenderink, G., and Siebesma, A.: Technical report ; TR - 302. The KNMI regional atmospheric climate model RACMO version 2.1, [Available online at <http://bibliotheek.knmi.nl/knmipubTR/TR302.pdf>], cited 2020.
- 705 Wadey, M. P., Haigh, I., Nicholls, R. J., Brown, J. M., Horsburgh, K., Carroll, B., Gallop, S. L., Mason, T., Bradshaw, E., et al.: A comparison of the 31 January–1 February 1953 and 5–6 December 2013 coastal flood events around the UK, *Frontiers in Marine Science*, 2, 84, 2015.
- Walters, D., Baran, A. J., Boutle, I., Brooks, M., Earnshaw, P., Edwards, J., Furtado, K., Hill, P., Lock, A., Manners, J., et al.: The Met Office Unified Model global atmosphere 7.0/7.1 and JULES global land 7.0 configurations, *Geoscientific Model Development*, 12, 1909–1963, 2019.
- 710 Weibull, W.: A statistical theory of strength of materials, IVB-Handl., 1939.
- Williams, J., Horsburgh, K. J., Williams, J. A., and Proctor, R. N.: Tide and skew surge independence: New insights for flood risk, *Geophysical Research Letters*, 43, 6410–6417, 2016.
- Williams, K., Copsey, D., Blockley, E., Bodas-Salcedo, A., Calvert, D., Comer, R., Davis, P., Graham, T., Hewitt, H., Hill, R., et al.: The Met Office global coupled model 3.0 and 3.1 (GC3. 0 and GC3. 1) configurations, *Journal of Advances in Modeling Earth Systems*, 10, 357–380, 2018.
- 715 Williams, K. D., Harris, C. M., Bodas-Salcedo, A., Camp, J., Comer, R. E., Copsey, D., Fereday, D., Graham, T., Hill, R., Hinton, T., Hyder, P., Ineson, S., Masato, G., Milton, S. F., Roberts, M. J., Rowell, D. P., Sanchez, C., Shelly, A., Sinha, B., Walters, D. N., West, A., Woollings, T., and Xavier, P. K.: The Met Office Global Coupled model 2.0 (GC2) configuration, *Geoscientific Model Development*, 8, 1509–1524, <https://doi.org/10.5194/gmd-8-1509-2015>, 2015.
- 720 Woollings, T., Hannachi, A., and Hoskins, B.: Variability of the North Atlantic eddy-driven jet stream, *Quarterly Journal of the Royal Meteorological Society*, 136, 856–868, 2010.