



# Machine-guided discovery of a real-world rogue wave model

Dion Häfner<sup>a,b,1</sup> , Johannes Gemmrich<sup>c</sup> , and Markus Jochum<sup>b</sup>

Edited by David Donoho, Stanford University, Stanford, CA; received April 19, 2023; accepted September 12, 2023

Big data and large-scale machine learning have had a profound impact on science and engineering, particularly in fields focused on forecasting and prediction. Yet, it is still not clear how we can use the superior pattern-matching abilities of machine learning models for scientific discovery. This is because the goals of machine learning and science are generally not aligned. In addition to being accurate, scientific theories must also be causally consistent with the underlying physical process and allow for human analysis, reasoning, and manipulation to advance the field. In this paper, we present a case study on discovering a symbolic model for oceanic rogue waves from data using causal analysis, deep learning, parsimony-guided model selection, and symbolic regression. We train an artificial neural network on causal features from an extensive dataset of observations from wave buoys, while selecting for predictive performance and causal invariance. We apply symbolic regression to distill this black-box model into a mathematical equation that retains the neural network's predictive capabilities, while allowing for interpretation in the context of existing wave theory. The resulting model reproduces known behavior, generates well-calibrated probabilities, and achieves better predictive scores on unseen data than current theory. This showcases how machine learning can facilitate inductive scientific discovery and paves the way for more accurate rogue wave forecasting.

ocean waves | rogue waves | machine learning | symbolic regression | causality

Rogue waves are extreme ocean waves that have caused countless accidents, often with fatal consequences (1). They are defined as waves whose crest-to-trough height  $H$  exceeds a threshold relative to the significant wave height  $H_s$ . The significant wave height is defined as four times the SD of the sea surface elevation. Here, we use a rogue wave criterion with a threshold of 2.0:

$$H/H_s > 2.0. \quad [1]$$

A rogue wave is therefore by definition an unlikely sample from the tail of the wave height distribution and can in principle occur by chance under any circumstance. This makes them difficult to analyze and requires massive amounts of data. Therefore, research has mostly focused on theory and idealized experiments in wave tanks, often considering only 1-dimensional wave propagation (2). However, the availability of large observation arrays (3) makes them an ideal target for machine-learning based analysis (4, 5).

In this study, we present a neural network-based model that predicts rogue wave probabilities from the sea state, trained solely on observations from buoys (6). The resulting model respects the causal structure of rogue wave generation; therefore, it can generalize to unseen physical regimes, is robust to distributional shift, and can be used to infer the relative importance of rogue wave generation mechanisms.

While a causally consistent neural network is useful for prediction and qualitative insight into the physical dynamics, the ability for scientists to analyze, test, and manipulate a model is crucial to recognize its limitations and integrate it into the research canon. Despite advances in interpretable AI (7), this is still a major challenge for most machine learning models.

To address this, we transform our neural network into a concise equation using symbolic regression (8, 9). The resulting model combines several known wave dynamics, outperforms current theory in predicting rogue wave occurrences, and can be interpreted within the context of wave theory. We see this as an example of “data-mining inspired induction” (10), an extension to the scientific method in which machine learning guides the discovery of new scientific theories.

We achieve this through the following recipe (Fig. 1):

1. A priori analysis of causal pathways that leads to a set of presumed causal parameters (Section 1).

## Significance

Machine learning has had a transformative impact on predictive science and engineering. But due to their black-box nature, better machine learning models do not always lead to greater human understanding, the first goal of science. We show how this can be overcome by using machine learning to transform a vast database of wave observations into a human-readable equation for the occurrence probability of rogue waves—rare ocean waves that routinely damage ships and offshore structures. This equation can be analyzed and incorporated into the research canon. Our work demonstrates the potential of causal analysis, machine learning, and symbolic regression to drive scientific discovery in a real-world application.

Author affiliations: <sup>a</sup>Pasteur Labs, Brooklyn, NY 11205; <sup>b</sup>Niels Bohr Institute, University of Copenhagen, Copenhagen 2100, Denmark; and <sup>c</sup>Department of Physics and Astronomy, University of Victoria, Victoria, BC V8W 2Y2, Canada

Author contributions: D.H., J.G., and M.J. designed research; D.H. performed research; D.H., J.G., and M.J. analyzed data; and D.H. wrote the paper.

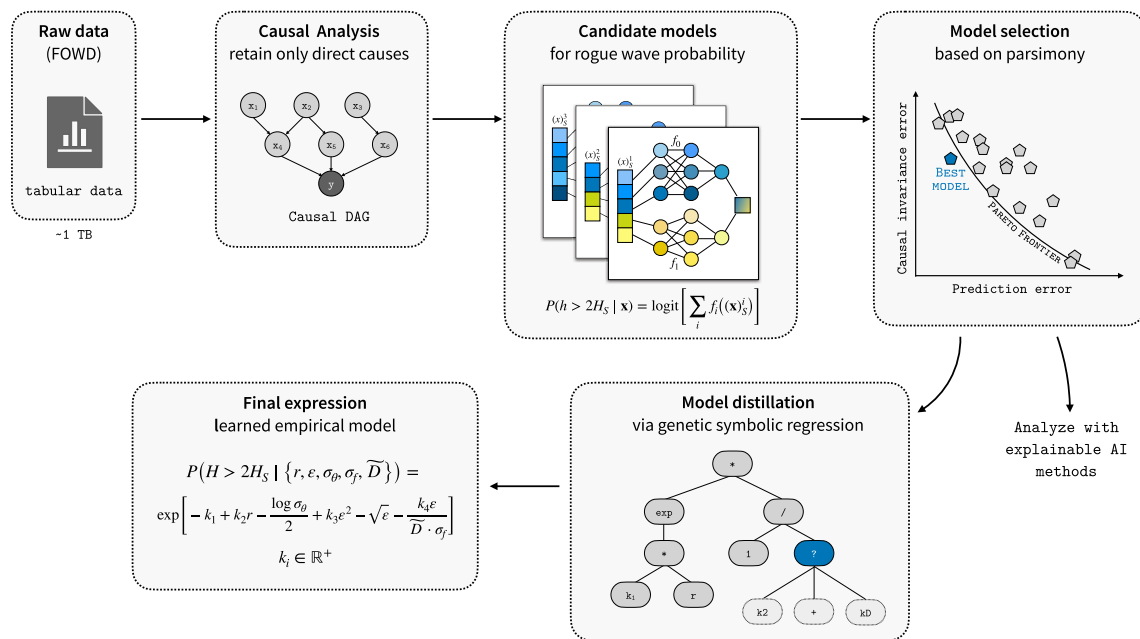
The authors declare no competing interest.

This article is a PNAS Direct Submission.

Copyright © 2023 the Author(s). Published by PNAS. This article is distributed under [Creative Commons Attribution-NonCommercial-NoDerivatives License 4.0 \(CC BY-NC-ND\)](#).

<sup>1</sup>To whom correspondence may be addressed. Email: [dion.haefner@simulation.science](mailto:dion.haefner@simulation.science).

Published November 20, 2023.



**Fig. 1.** Overview of our study. Starting out with large amounts of tabular data from wave buoys, we use a causal analysis to identify the most important features for predicting rogue waves. We then train an ensemble of neural networks on subsets of these features and select the best one based on its predictive performance and causal invariance. Finally, we use symbolic regression to distill the model into a concise mathematical equation. We analyze the neural network and symbolic expression in terms of their performance on unseen data and compare them to existing theory. This closes the arc between data, machine learning, and theory.

2. Training an ensemble of regularized neural network predictors, and parsimony-guided model selection based on causal invariance (Section 2).
3. Distillation of the neural network into a concise mathematical expression via symbolic regression (Section 3).

Finally, we analyze both the neural network and symbolic model in the context of current wave theory (Section 4). Both models reproduce well-known behavior and point toward insights regarding the relative importance of different mechanisms in the real ocean.

## 1. A Causal Graph for Rogue Wave Generation

To create a causal machine learning model, it is crucial to expose it only to parameters with causal relevance. Otherwise, the model may prefer to encode spurious associations over true causal relationships, simply because they can be easier to learn. This requires us to identify the causal structure of rogue wave generation.

There are several hypothesized causes of rogue waves see ref. 11, for an overview. Typically, research focuses on linear superposition in finite-bandwidth seas (12), wave breaking (13), and wave-wave interactions in weakly nonlinear seas (14, 15) or through the modulational instability (16). Apart from these universal mechanisms, there are also countless possible interactions with localized features such as nonuniform topography (17), wave-current interactions like in the Agulhas (18) or the Antarctic Circumpolar Current (19), or crossing sea states at high crossing angles affecting wave breaking (20). We call this set of mechanisms the physical effects  $\Phi$ .

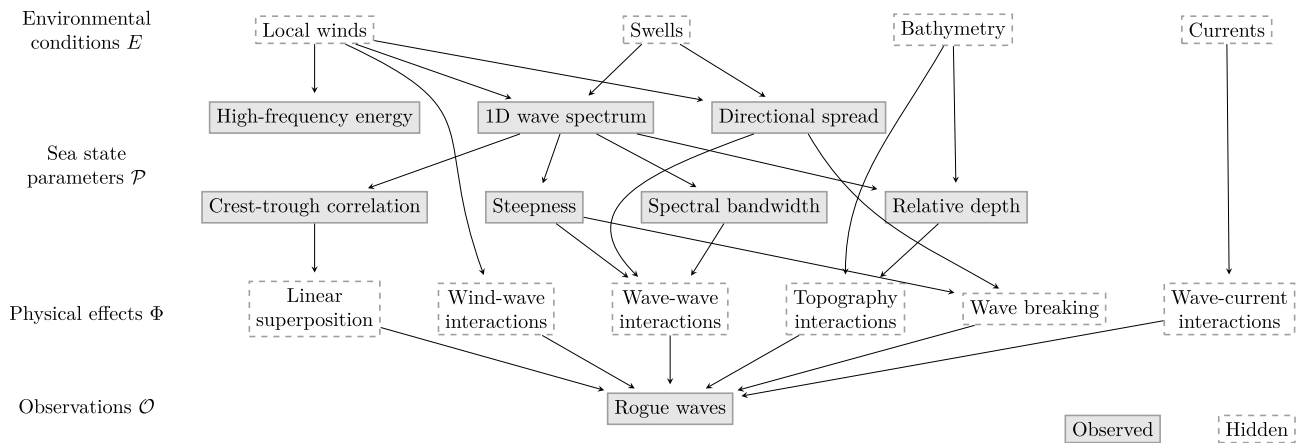
Since ocean waves are generated by a complex dynamical system, their true cause is a set of extrinsic environmental conditions  $E$  that are high-dimensional and not feasible to capture in full detail. However, most physical effects are mediated by

one or several sea state parameters  $\mathcal{P}$ , which are the characteristic aggregated parameters that appear in theoretical models of the respective wave dynamics and that are included in operational wave forecasts. In this study, we would like to obtain a model that relates relevant sea state parameters  $\mathcal{P}$  to wave observations  $\mathcal{O}$ , which ideally also lets us infer the relative importance of physical effects  $\Phi$ .

The go-to tool to analyze causal relationships is a causal DAG [Directed Acyclic Graph; (21)]. In a causal DAG, nodes represent variables and edges  $A \rightarrow B$  imply that  $A$  is a cause of  $B$  [usually in the probabilistic sense in that the probability distribution  $P(B)$  depends on  $A$ ].

We create a causal graph for rogue wave formation based on the hypothesized causal mechanisms discussed above and their corresponding theoretical models and parameters (Fig. 2). Following this causal structure, we use the following set of sea state parameters as candidates for representing the various causal pathways (see the *Materials and Methods* for more information on each parameter):

- Crest–trough correlation  $r$ , to account for the linear effect of wave groups on crest-to-trough rogue waves (22).  $r$  is the dominant causal factor behind linear rogue wave formation (4).
- Steepness  $\epsilon$  governing weakly nonlinear effects, such as second-order and third-order bound waves, and wave breaking (13, 23).
- Relative high-frequency energy  $E_h$  (fraction of total energy contained in the spectral band 0.25 Hz to 1.5 Hz) as a proxy for the strength of local winds (24).
- Relative depth  $\tilde{D}$  (based on peak wavelength), which is central for nonlinear shallow-water effects (25, 26) and wave breaking (13).
- Dominant directional spread  $\sigma_\theta$ , which has an influence on third-order nonlinear waves (26) and wave breaking (20).



**Fig. 2.** The causes of rogue waves as a causal DAG (directed acyclic graph). Arrows  $A \rightarrow B$  imply that  $A$  causes  $B$ .

- Spectral bandwidth  $\nu_f$  (narrowness) and  $\sigma_f$  (peakedness), appearing, for example, in the expression for the influence of third-order nonlinear waves (26).

We also include a number of derived parameters that commonly appear in wave models and govern certain nonlinear (wave-wave) phenomena:

- Benjamin–Feir index BFI, which controls third-order nonlinear free waves (26) and the modulational instability (27).
- Ursell number  $Ur$ , which quantifies nonlinear effects in shallow water (28).
- Directionality index  $R$  (the ratio of directional spread and spectral bandwidth), which has an influence on third-order nonlinear free waves and is typically used in conjunction with the BFI (26).

These parameters cover most causal pathways toward rogue wave generation. Still, there are some at least partially unobserved causes, as we do not have access to data on local winds, topography, or currents. Additionally, our in situ measurements are potentially biased estimates of the true sea state parameters, and there is no guarantee that any given training procedure will converge to the true causal model. This implies that we cannot rely on a model being causally consistent by design; instead, we perform a posteriori verification on the learned models to find the perfect trade-off between causal consistency and predictive performance (Section 2C).

## 2. An Approximately Causal Neural Network

**A. Input Data.** We use the Free Ocean Wave Dataset [FOWD, (6)], which contains 1.4 billion wave measurements recorded by the 158 CDIP wave buoys (3) along the Pacific and Atlantic coasts of the US, Hawaii, and overseas US territories. Water depths range between 10 m to 4,000 m, and we require a significant wave height of at least 1 m. Each buoy records the sea surface elevation at a sampling frequency of 1.28 Hz, producing over 700 y of time series in total. FOWD extracts every zero-crossing wave from the surface elevation data and computes a number of characteristic sea state parameters from the history of the wave within a sliding window.

Due to the massive data volume of the full FOWD catalogue ( $\sim 1$  TB), we use an aggregated version that maps each sea state to the maximum wave height of the following 100 waves as in

ref. 4. This reduces the data volume by a factor of 100 and inflates all rogue wave probabilities to a bigger value  $\hat{p}$ . We correct for this via  $p = 1 - (1 - \hat{p})^{1/100}$ , assuming that rogue waves occur independently from each other. This is a good approximation in most conditions but may underestimate seas with a strong group structure (Section 5B).

The final dataset has 12.9M data points containing over 100,000 rogue waves exceeding 2 times the significant wave height. Our dataset is freely available for download (Data, Materials, and Software Availability).

**B. Neural Network Architecture.** The probability to measure a rogue wave based on the sea state can be modeled as a sum of nonlinear functions, each of which only depends on a subset of the sea state parameters representing a different causal path (act via different physical effects in Fig. 2):

$$\text{logit } P(y = 1 \mid \mathbf{x}) \sim \sum_i f_i(\mathbf{x}^{(S_i)}) + b. \quad [2]$$

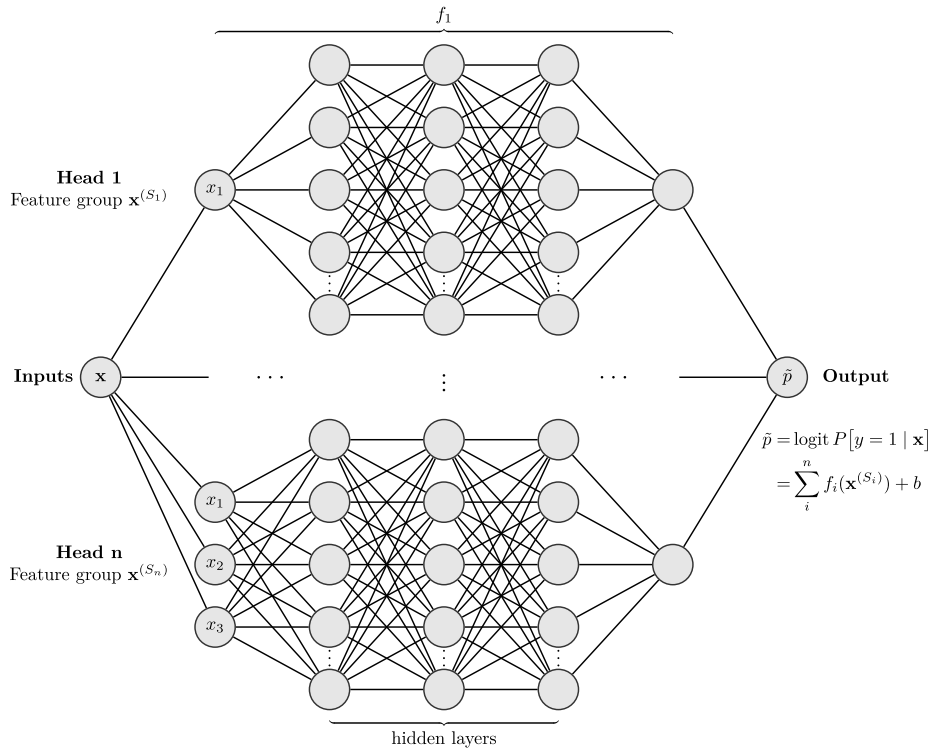
Here,  $y$  is a binary label indicating whether the current wave is a rogue wave,  $\mathbf{x}^{(S_i)}$  is the  $i$ -th subset of all causal sea state parameters  $\mathbf{x}$ ,  $\text{logit}(p) = \log(p) - \log(1 - p)$  is the logit function,  $f_i$  are arbitrary nonlinear functions to be learned, and  $b$  is a constant bias term.

By including only a subset  $\mathbf{x}^{(S_i)}$  of all parameters  $\mathbf{x}$  as input for  $f_i$ , we can restrict which parameters may interact nonadditively with each other, which is an additional regularizing constraint that increases interpretability and prevents interactions between inputs from different causal pathways. For example, to include the effects of linear superposition and nonlinear corrections for free and bound waves as in ref. 29, Eq. 2 can be written as

$$\text{logit } P(y = 1 \mid \mathbf{x}) \sim \underbrace{f_1(r)}_{\text{linear}} + \underbrace{f_2(\text{BFI}, R)}_{\text{free waves}} + \underbrace{f_3(\epsilon, \tilde{D})}_{\text{bound waves}}. \quad [3]$$

We use a neural network with fully connected layers (FCN) to model the functions  $f_i$ , which are universal function approximators (30), and that can be trained efficiently for large amounts of data. The set of functions  $f_i$  can be represented as a single multihead FCN with a linear output layer (Fig. 3). We use a small feed-forward architecture with 3 hidden layers and ReLU activation functions [rectified linear units, (31)].

The neural network outputs a scalar  $\tilde{p} \in (-\infty, \infty)$ , the log-odds of a rogue wave occurrence for the given sea state. For



**Fig. 3.** Neural network architecture (multihead FCN) used to predict rogue wave probabilities. Each input head receives a different subset of the full parameter set  $\mathbf{x}$  to limit the amount of noncausal interactions between parameters.

training, we use the Adam optimizer (32) and backpropagation to minimize a cross-entropy loss for binary classification with an added  $\ell_2$  regularization term for kernel parameters:

$$L(p, y, \theta) = y \cdot \log(p) + (1 - y) \cdot \log(1 - p) + \lambda \|\theta\|_2, \quad [4]$$

with predicted probability  $p = \text{logit}^{-1}(\tilde{p})$ , observed labels  $y \in \{0, 1\}$  (rogue wave or not), and neural network kernel parameters  $\theta$ .

To estimate uncertainties in the neural network parameters and resulting predictions, we use Gaussian stochastic weight averaging [SWAG, (33)]. For this, we train the network for 50 epochs, then start recording the optimizer trajectory after each epoch for another 50 epochs. The observed covariance structure of the sampled parameters is used to construct a multivariate Gaussian approximation of the loss surface that we can sample from. This results in slightly better predictions and gives us a way to quantify how confident the neural network is in its predictions.

**C. Causal Consistency and Predictive Accuracy.** Although we include only input parameters that we assume to have a direct causal connection with rogue wave generation, there is no guarantee that the neural network will infer the correct causal model. In fact, the presence of measurement bias and unobserved causal paths makes it unlikely that the model will converge to the true causal structure. To search for an approximately causally consistent model, we will have to quantify its causal performance.

We achieve this through the concept of invariant causal prediction [ICP; (34, 35)]. The key insight behind ICP is that the parameters of the true causal model will be invariant under distributional shift, that is, an intervention on an upstream “environment” node in the causal graph that controls which distribution the data are drawn from. Retraining the model on data with different spurious correlations between features should

still lead to the same dependency of the target on the features see also ref. 36.

We split the dataset randomly into separate training and validation sets, in chunks of 1M waves. We train the model on the full training dataset and perform ICP on the validation dataset, which we partition into subsets representing different conditions in space, time, depth, spectral properties, and degrees of nonlinearity (Table 1). This changes the dominant characteristics of the waves in each subset (representing, e.g. storm and swell conditions), inducing distributional shift. Then, we retrain the model separately on each subset and compute the Rms difference between predictions of the retrained model  $P_k$  and the full model  $P_{\text{tot}}$  on the  $k$ -th data subset  $\mathbf{x}_{(k)}$ :

$$\mathcal{E}_k^2 = \frac{1}{n_k} \sum_i^{n_k} \left( \text{logit } P_k(\mathbf{x}_i^{(k)}) - \text{logit } P_{\text{tot}}(\mathbf{x}_i^{(k)}) \right)^2. \quad [5]$$

As the total consistency error, we use the rms of Eq. 5 across all environments:

$$\mathcal{E} = \sqrt{\frac{1}{n_E} \sum_k^{n_E} \mathcal{E}_k^2}. \quad [6]$$

Under a noise-free, infinite dataset and an unbiased training process that always identifies the true causal model we would find  $\mathcal{E} = 0$ , i.e., retraining the model on the unseen data subset would not contribute any new information and leave the model perfectly invariant. Since all of these assumptions are violated here, we merely search for an approximately causal model that minimizes  $\mathcal{E}$ .

However, we cannot use  $\mathcal{E}$  as the only criterion when selecting a model. The invariance error can only account for change in the prediction (variance), but not for its overall closeness to the

**Table 1. The subsets of the validation dataset used to evaluate model performance and invariance**

Subset name	Condition	# waves
Southern-California	Longitude $\in (-123.5, -117)^\circ$ , latitude $\in (32, 38)^\circ$	265M
Deep-stations	Water depth $> 1,000$ m	28M
Shallow-stations	Water depth $< 100$ m	154M
Summer	Day of year $\in (160, 220)$	51M
Winter	Day of year $\in (0, 60)$	91M
$H_s > 3$ m	$H_s > 3$ m	58M
High-frequency	Relative swell energy $< 0.15$	43M
Low-frequency	Relative swell energy $> 0.7$	46M
Long-period	Mean zero-crossing period $> 9$ s	100M
Short-period	Mean zero-crossing period $< 6$ s	42M
Cnoidal	Ursell number $> 8$	40M
Weakly-nonlinear	Steepness $> 0.04$	83M
Low-spread	Directional spread $< 20^\circ$	25M
High-spread	Directional spread $> 40^\circ$	25M
Full	(all validation data)	472M

true solution (bias). Therefore, we select a model that is Pareto-optimal with respect to the invariance error  $\mathcal{E}$  and a predictive score  $\mathcal{L}$ . This will not establish absolute causal consistency but will allow us to select a model that is near-optimal given the constraints.

For  $\mathcal{L}$ , we use the log of the likelihood ratio between the predictions of our neural network and a baseline model that predicts the empirical base rate  $\bar{y}_k = \frac{1}{n} \sum_i^n y_{k,i}$ , averaged over all environments  $k$ :

$$\mathcal{L}(p, \bar{y}) = \frac{1}{n_E} \sum_k (I(p_k) - I(\bar{y}_k)), \quad [7]$$

$$I(x) = x \cdot \log(x) + (1 - x) \cdot \log(1 - x). \quad [8]$$

To evaluate model calibration (the tendency to produce over- or underconfident probabilities), we compute a calibration curve by binning the predicted rogue wave probabilities. We then compare each bin to the observed rogue wave frequency, and compute the weighted rms residual between measured ( $\bar{y}_i$ ) and predicted ( $p_i$ ) log-odds:

$$\mathcal{C} = \sqrt{\sum_{i=1}^{n_b} w_i (\text{logit}(p_i) - \text{logit}(\bar{y}_i))^2}. \quad [9]$$

To account for uncertainty in the observations (e.g., close to the extremes), the weights  $w_k$  are based on the 33% credible interval of  $\bar{y}_i \sim \text{Beta}(n_i^+, n_i^-)$  with  $n_i^+$  rogue and  $n_i^-$  nonrogue measurements. This is similar to the expected calibration error (37) but models data uncertainty directly. We use a uniform bin size (in logit space) of 0.1.

**D. Model Selection.** We train a total of 24 candidate models on different subsets of the relevant causal parameters (as identified in Section 1) and varying number of input heads (between 1 and 3). We evaluate their performance in terms of calibration, predictive performance, and causal consistency (Table 3).

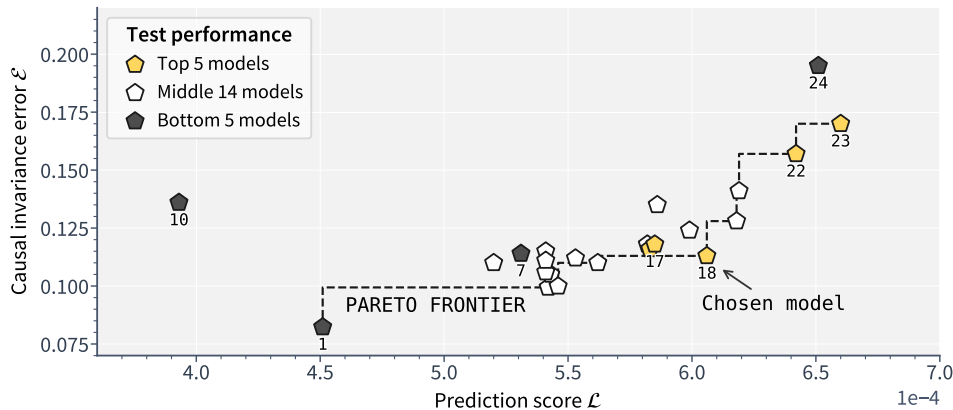
We observe a clear anticorrelation between model complexity and predictive score on one hand and causal consistency on the other hand (Fig. 4). This is evidence that more complex models are indeed less biased but exploit more noncausal connections. We perform model selection based on parsimony: A good model is one where a small increase in either predictive performance or causal consistency implies a large decrease in the other, i.e., where the Pareto front is convex. This is similar to the metric used by PySR (9) to select the best symbolic regression model (Section 3).

Based on this, we choose model 18 with parameter groups  $S_1 = \{r\}$ ,  $S_2 = \{\epsilon, \sigma_\theta, \sigma_f, \tilde{D}\}$  (i.e., a model with two input heads) as the reference model for further analysis. The chosen model produces well-calibrated probabilities (Fig. 5) and is among the 5 best models in terms of predictive performance on the test dataset (not used during training or selection), despite using only 5 features with at most 4-way interactions.

The relatively low number of input features allows us to analyze the model in detail using explainable AI methods (Section 4A).

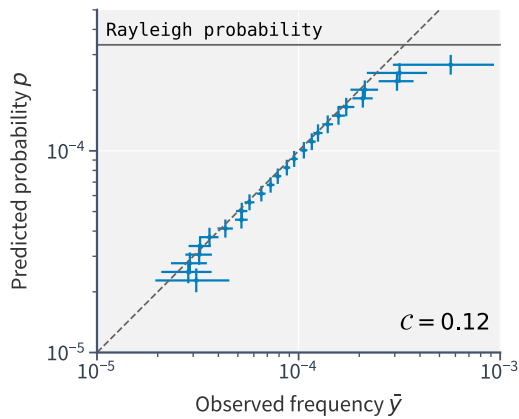
### 3. Learning an Empirical Equation for Rogue Wave Risk

To make our model fully interpretable, we transform the learned neural network into an equation via symbolic regression. Common approaches to symbolic regression include Eureqa (39), AI Feynman (40), SINDy (41), and QLattice (42). Here,



**Fig. 4.** There is a clear trade-off between causal invariance ( $\mathcal{E}$ ) and predictive performance ( $\mathcal{L}$ ) of our neural network predictors. We choose the model that lies in the most convex part of the Pareto frontier. Scores are evaluated on validation data. Test performance is based on prediction scores on held-out test data (from unseen stations).





**Fig. 5.** Our model outputs well-calibrated probabilities, even for unseen stations. Shown is the binned predicted probability  $p$  vs. the observed rogue wave frequency  $\bar{y}$  on the test data. Error bars for  $p$  indicate 3 standard deviations estimated via SWAG sampling. Error bars for  $\bar{y}$  indicate 95% credible interval assuming  $\bar{y}_i \sim \text{Beta}(\eta_i^+, \eta_i^-)$ . Bins with less than 10 observed rogue waves are excluded. The dashed line indicates perfect calibration. Solid line indicates probability as predicted by linear theory in the narrow-bandwidth limit [Rayleigh distribution; (38)].

we use PySR (8, 9), a symbolic regression package based on genetic programming (43). Genetic algorithms build a large ensemble of candidate models and select the best ones, before mutating and recombining them into the next generation. In the case of symbolic regression, mathematical expressions are represented as a tree of constants and elementary symbols. In principle, this allows PySR to discover expressions of unbounded complexity.

PySR's central metric to quantify the goodness of an equation is again based on parsimony, in the form of the derivative of predictive performance with respect to the model complexity—if the true model has been discovered, any additional complexity can at best lead to minor performance gains (by overfitting to noise in the data).

In our case, we seek to find an expression  $f$  from the space of possible expression graphs  $\mathcal{T}_O$  with allowed operators  $O$  that approximates the rogue wave log-probability as predicted by the neural network  $\mathcal{N}$  over the dataset  $x$ :

$$\text{Find } f \in \mathcal{T}_O \text{ that minimizes } \sum_i \frac{1}{\text{Var}(y_i)} \left[ f(x_i)^2 - \sigma(\mathbb{E}[y_i])^2 \right],$$

where  $\sigma(x) = -\log(1 + \exp(-x))$ , and  $y_i$  is the set of SWAG samples from  $\mathcal{N}(x_i)$ . A sensible set of operators  $O$  is key to ensure interpretability of the resulting expression; we choose the symbols  $O = \{+, -, \times, \div, \log, \cdot^{-1}, \sqrt{\cdot}, \cdot^2\}$  to facilitate expressions that are similar to current theoretical models of the form  $P \sim A \exp(B)$ . We normalize all input features to approximately unit scale by converting directional spread to radians.

PySR assembles a league of candidate expression and presents the Pareto-optimal solutions of increasing complexity to the user. We select the best solution by hand, picking the expression with the best parsimony score that contains all input features and at least two terms containing the steepness  $\varepsilon$  (to account for the various causal pathways in which steepness affects rogue waves). The final equation is shown in Fig. 8, and discussed in Section 4B.

## 4. Results

**A. Neural Network.** We analyze the behavior of our neural network predictor, which reveals important insights about the physical dynamics of rogue waves and their prediction.

**A.1. Rogue wave models should account for crest-trough correlation, steepness, relative depth, and directionality.** Only this parameter combination achieves good causal consistency and predictive scores at the same time, and experiments that exclude any of these parameters perform unconditionally worse in either metric. Especially the exclusion of crest–trough correlation leads to catastrophic results, even when including other bandwidth measures like  $\sigma_\theta$  in its place (Table 3).

This suggests that the above set of parameters represents the dominant rogue wave generation processes in the form of linear superposition in finite-bandwidth seas with a directional contribution and weakly nonlinear corrections.

The crest–trough correlation  $r$  is still lacking mainstream adoption as a rogue wave indicator for example, it is not part of ECMWF's operational forecast (29), despite being a key parameter for crest-to-trough rogue waves (4, 22, 44). The other parameters are consistent with other empirical studies such as Fedele (45), which considers the same parameters in conjunction with rogue crests during storms. They are also similar to the ingredients to ECMWF's rogue wave forecast (29), which is based on the effects of second and third-order bound and free waves and uses steepness, relative depth, directional spread, and spectral bandwidth. However, in our model, these parameters are combined differently; a model enforcing the same interactions (steepness and relative depth for bound wave contribution, BFI, and directionality index for free wave contribution) performs poorly.

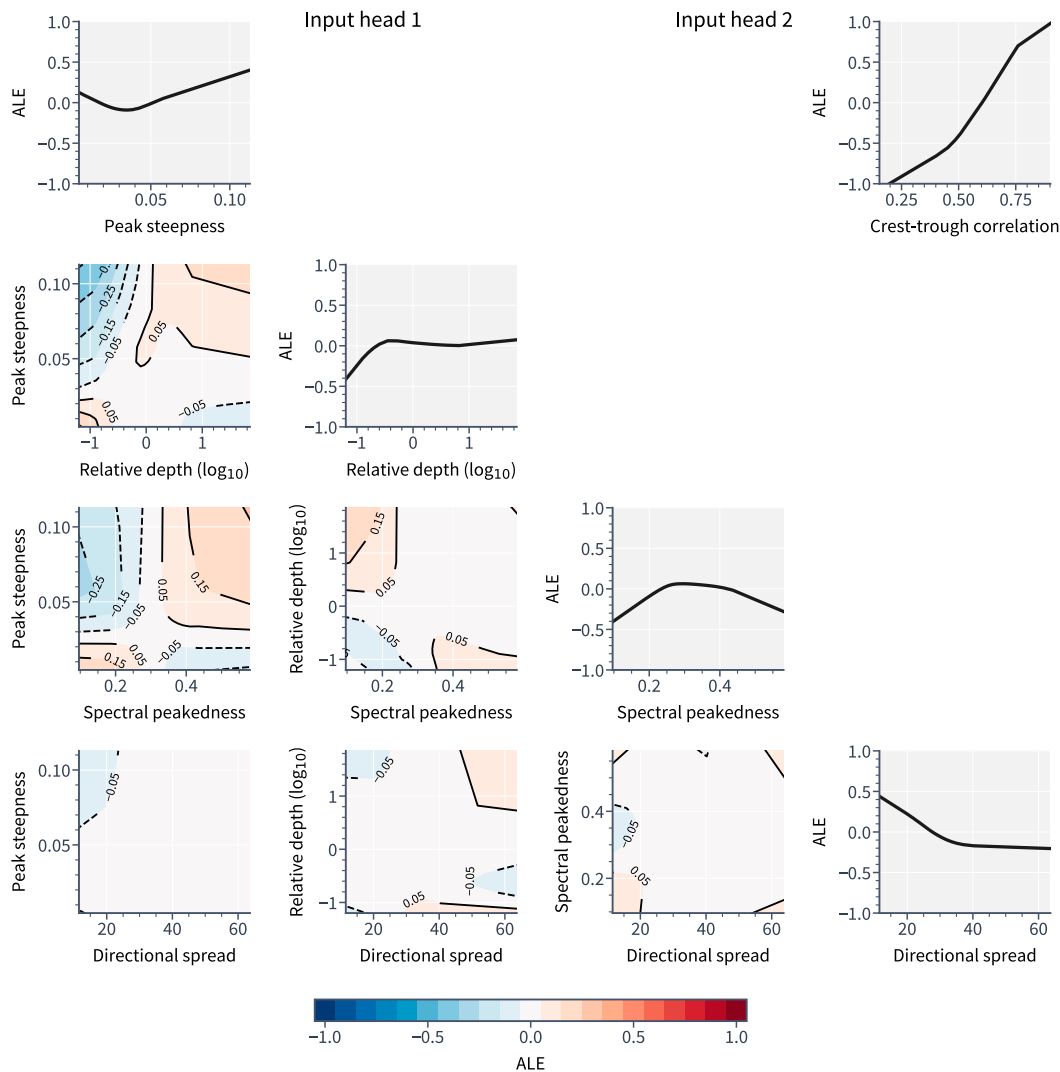
Numerous previous studies have found the BFI to be a poor predictor of rogue wave risk in realistic sea states (4, 14, 15, 45–48) due to its strong underlying assumptions such as unidirectionality. This study extends this to the fully nonparametric and nonlinear case.

We study how our model uses different parameters by visualizing their impact on the prediction of the respective head of the neural network. For this, we make use of the accumulated local effects decomposition [ALE, (49)], which measures the influence of infinitesimal changes in each parameter on the prediction outcome see also ref. 7. From the ALE plot (Fig. 6), we find that crest–trough correlation has by far the biggest influence of all parameters and explains about 1 order of magnitude in rogue wave risk variation, which is consistent with earlier model-free approaches (4). To first order, higher crest–trough correlation, lower directional spread, larger relative depth (deep water), and higher steepness lead to larger rogue wave risk, but parameter interactions can lead to more complicated, nonmonotonic relationships (for example, in very shallow water; see Section 4A.3).

**A.2. The Rayleigh distribution is an upper bound for real-world rogue wave risk.** Despite the clear enhancement by weakly nonlinear corrections, the Rayleigh wave height distribution remains an upper bound for real-world (crest-to-trough) rogue waves. The Rayleigh distribution is the theoretical wave height distribution for linear narrow-band waves (38), i.e., the limit  $r \rightarrow 1$ ,  $\varepsilon \rightarrow 0$ ,  $\sigma_f \rightarrow 0$ ,  $\tilde{D} \rightarrow \infty$ , and  $\sigma_\theta \rightarrow 0$ , and reads:

$$P(H/H_s > k) = \exp(-2k^2). \quad [10]$$

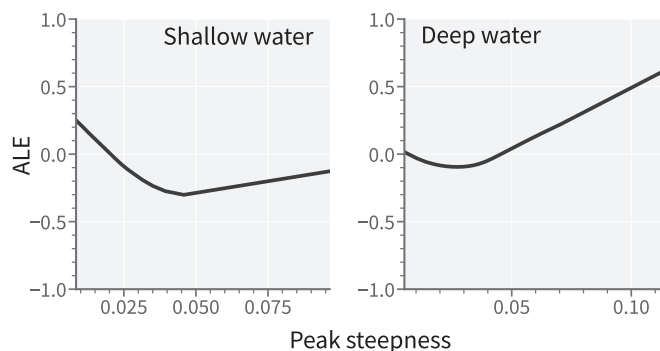
Only in the most extreme conditions does our model predict a similarly high probability, for example for  $\sigma_\theta = 13^\circ$ ,  $\varepsilon = 0.008$ ,



**Fig. 6.** ALE (accumulated local effects) plot matrix for experiment 18. Shown is the change in rogue wave risk (in logits) from the average as each parameter is varied. The total effect is the sum of all 1D, 2D, and higher-order contributions (not shown).

$\sigma_f = 0.14$ ,  $r = 0.88$ , and  $\tilde{D} = 0.6$ , which gives the same probability as the Rayleigh distribution,  $p = 3.3 \times 10^{-4}$ .

In the opposite extreme, rogue wave probabilities can fall to as little as  $1 \times 10^{-5}$  for low values of  $r$  and high values of  $\sigma_\theta$



**Fig. 7.** Our model predicts a positive association between steepness and rogue waves in deep water and a negative association in shallow water. Shown is the 1-dimensional ALE (accumulated local effects) plot in both cases. Here, deep water corresponds to sea states with  $\tilde{D} > 3$  and shallow water with  $\tilde{D} < 0.1$ .

(such as in a sea with a strong high-frequency component and high directional spread). This suggests that bandwidth effects can create sea states that efficiently suppress extremes.

**A.3. There is a clear separation between deep water and shallow water regimes.** All models with high causal invariance scores include an interaction between steepness and relative water depth. Looking at this more closely, we find that a stratification on deep and shallow water sea states reveals 2 distinct regimes (Fig. 7).

In deep water, rogue wave risk is strongly positively associated with steepness, as expected from the contribution of second and third-order nonlinear bound waves (26). The opposite is true in shallow water ( $\tilde{D} < 0.1$ ), where we find a clear negative association with steepness. This is likely due to depth-induced wave breaking (23). In very shallow waters, more sea states have a steepness close to the breaking threshold, which removes taller waves that tend to have a higher steepness than average (Fig. 8).

**B. Symbolic Expression.** The final expression for the rogue wave probability, as discovered via symbolic regression, is given in Fig. 8. It consists of an exponential containing five additive terms:

$$P(H > 2H_S | r, \varepsilon, \sigma_\theta, \sigma_f, \tilde{D}) = \exp \left[ \underbrace{-12. + 3.8r}_{\text{I}} - \underbrace{\frac{\log \sigma_\theta}{2}}_{\text{II}} + \underbrace{66.\varepsilon^2}_{\text{III}} - \underbrace{\sqrt{\varepsilon}}_{\text{IV}} - \underbrace{\frac{0.23\varepsilon}{\tilde{D} \cdot \sigma_f}}_{\text{V}} \right]$$

**Fig. 8.** Our empirical equation for rogue wave risk, as identified through the distillation of our neural network predictor via symbolic regression. This equation outperforms existing wave theory on unseen stations from our dataset, while being fully interpretable. Numbered terms are discussed in Section 4B. All floating point coefficients are rounded to two significant digits.

- (I)  $-12 + 3.8r$ . The term with the largest coefficients is the one containing  $r$ , as expected. Comparison with the exponential term in the Tayfun distribution  $P_t$ , Eq. 28, reveals that this is approximately a linear expansion around  $r \approx 1$ :

$$\log P_t(H/H_s > h) \sim -\frac{4b^2}{1+r}, \quad [11]$$

$$= -12 + 4r + \mathcal{O}(r^2) \Big|_{r \approx 1}. \quad [12]$$

This is an important sanity check for the model since it shows that it is able to rediscover existing theory purely from data.

- (II)  $-\log \sigma_\theta/2$ . This encodes the observed enhancement for narrow sea states and has no direct relation to existing quantitative theory. Its functional form is somewhat problematic since it causes the model to diverge for  $\sigma_\theta \rightarrow 0$  (unidirectional seas). However, the model has only seen real-world seas with  $\sigma_\theta \gtrsim 0.2$ , so we may replace this term with one that yields similar predictions for the relevant range of  $\sigma_\theta$ , and does not diverge for  $\sigma_\theta \rightarrow 0$ .

One possible candidate is

$$\frac{1 - \sigma_\theta}{1 + \sigma_\theta}, \quad [13]$$

which has a relative RMS error of about 5% over the range  $\sigma_\theta \in (20, 90)^\circ$  compared to the original term.

- (III)  $66\varepsilon^2$ . Encodes the influence of weakly nonlinear effects for large values of  $\varepsilon \gtrsim 0.1$ .
- (IV)  $-\sqrt{\varepsilon}$ . This term encodes the observed negative association between steepness and rogue waves for low values of  $\varepsilon$  that could be due to wave breaking or may be an artifact of our sensor.

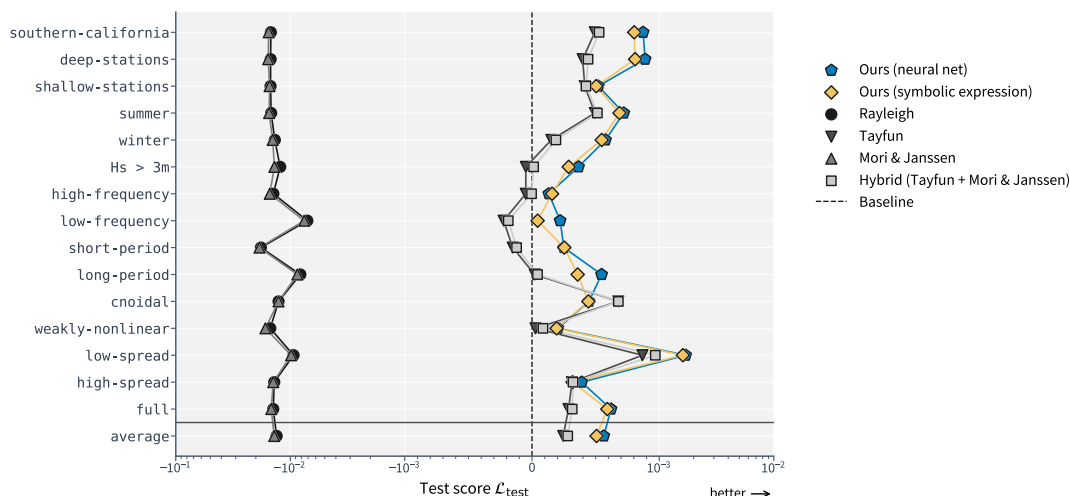
- (V)  $0.23\varepsilon/(\tilde{D} \cdot \sigma_f)$ . Since  $\tilde{D} \sim k_p D$  and  $\varepsilon \sim k_p H_s$ , this term is proportional to the relative wave height  $\eta = H_s/D$  and  $1/\sigma_f$ .  $\eta$  is the most important parameter in the theory of shallow-water waves and appears for example in the Korteweg-de Vries equation (25). Accordingly, this term dominates the dynamics in very shallow water. Dependencies on  $1/\sigma_f$  occur in current theory (26) but are usually paired with  $\sigma_\theta$  to form the directionality index  $R$ . This suggests that term V may be incomplete and missing physical dynamics that are not prevalent in the data.

Overall, the equation is able to reproduce the same qualitative behavior as observed from the neural network, with the same well-calibrated outputs ( $\mathcal{C} = 0.14$ ) and predictive performance (Section 5A) on the test data.

## 5. Discussion

**A. Validation Against Theory.** We test our models (neural network and symbolic equation) against existing wave theory based on their mean predictive score  $\mathcal{L}$  across the environments from Table 1 on the held-out test data (unseen stations). As theoretical baselines, we use the models from Longuet-Higgins (Rayleigh, 38), Tayfun (22), Mori & Janssen (50), and a hybrid combining Tayfun and Mori & Janssen (*Materials and Methods*).

The results are shown in Fig. 9. Since the Rayleigh and Mori & Janssen models do not account for crest–trough correlation, their predictions vastly overestimate the occurrence rate of observed rogue waves. The Tayfun and hybrid models perform better but are still outperformed by our models except in cnoidal seas. Our models are better predictors than the baseline (predicting the empirical per-environment rogue wave frequency) in all environments.



**Fig. 9.** Comparison between our models and existing theory on held-out test data. Our models perform similar to each other and outperform existing theory on this dataset in all but one data subset (cnoidal seas). x-scale is linear in  $(-1 \times 10^{-3}, 1 \times 10^{-3})$ , and logarithmic otherwise.



The neural network performs better than the symbolic equation in all environments, albeit only by a small margin. This shows that the symbolic equation is able to capture the main features of the full model, despite its compact representation.

**B. Limitations.** Using only wave buoy observations for our analysis, we acknowledge the following limitations:

- We did not have sufficient data on local winds, currents, or topography, which implies that some relevant causal pathways are unobserved (Fig. 2). While we expect these effects to play a minor role in bulk analysis, they could dramatically affect local rogue wave probabilities in specific conditions, for example, over sloping topography (17) or in strong currents (51).
- We only have one-dimensional (time series) data and cannot capture imported parameters, such as solitons generated elsewhere that travel into the observation area. While we expect this to play a minor role, it could underestimate the importance of nonlinear free waves.
- Systematic sensor bias is common in buoys and can lead to spurious causal relationships. This may obscure the true causal structure and hurt model generalization to other sensors. However, this adaptation to sensor characteristics may be desirable in forecasting scenarios, where it allows the model to synthesize several noisy quantities into more robust ones.
- By aggregating individual waves into 100-wave chunks, we underestimate the per-wave rogue wave probability in sea states in which rogue waves do not occur independently of each other, such as seas with a strong group structure.

These limitations could potentially reduce our model’s ability to detect relevant causal pathways and underestimate the true rogue wave risk. Our analysis is agnostic to the data source and can be repeated on different sources to validate our findings.

## 6. Next Steps

**A. An Improved Rogue Wave Forecast.** Our empirical model can be compared directly to existing rogue wave risk indicators by evaluating them on forecast sea state parameters. ECMWF’s operational rogue wave forecast (29) focuses on envelope wave heights which does not account for crest–trough correlation and is conceptually similar to the Mori & Janssen model in Section 5A. Therefore, we are confident that substantial improvements are within reach in terms of predicting crest-to-trough rogue waves, even without using a black-box model.

**B. Predicting Superrogue Waves.** Observed wave height distributions often show a flattening of the wave height distribution toward the extreme tail (11, 14, 52). Therefore, we expect rogue wave probabilities to be more pronounced for even more extreme waves for example with  $H/H_s > 2.4$ , as recently observed in ref. 53.

The lack of sufficient direct observations in these regimes calls for a different strategy. One approach could be to transform this classification problem (rogue wave or not) into a regression, where the predicted variables are the free parameters of a candidate wave height probability distribution (such as shape and scale parameters of a Weibull distribution). Then, a similar analysis as in this study could be conducted for these parameters, which may reveal the main mechanisms influencing the risk for truly exceptional waves, and whether this flattening can be confirmed in our dataset.

**C. Commoditization of Data-Mining Based Induction.** There is a pronounced lack of established methods for machine learning aimed at scientific discovery. We have shown that incorporating and enforcing causal structure can overcome many of the shortcomings of standard machine learning approaches, like poorly calibrated predictions, noninterpretability, and incompatibility with existing theory. However, the methods we leveraged are still in their infancy and rely on further community efforts to be end-to-end automated and adopted at scale. Particularly, parsimony-based model selection (as in Sections 2D and 3) is still a manual process that requires a firm understanding of model intrinsics and the domain at hand. Nonetheless, we believe that the potential benefits of causal and parsimony-guided machine learning for real-world problems are too great to ignore, and we hope that this study will inspire further research in this direction.

## Materials and Methods

**Sea State Parameters.** Here, we give the definition of the sea state parameters used in this study. For a more thorough description of how parameters are computed from buoy displacement time series, see Häfner et al. (6).

All parameters can be derived from the nondirectional wave spectrum  $S(f)$ , with the exception of directional spread  $\sigma_\theta$ , which is estimated from the horizontal motion of the buoy and taken from the raw CDIP data.

Most parameters are computed from moments of the wave spectrum, where the  $n$ -th moment  $m_n$  is defined as

$$m_n = \int_0^\infty f^n S(f) \, df. \tag{14}$$

The expressions for the relevant sea state parameters are as follows:

- Significant wave height:  $H_s = 4\sqrt{m_0},$  (15)

- Spectral bandwidth (narrowness):  $\nu_f = \sqrt{m_2 m_0 / m_1^2} - 1,$  (16)

- Spectral bandwidth (peakedness):  $\sigma_f = \frac{m_0^2}{2\sqrt{\pi}} \left( \int_0^\infty f \cdot S(f) \, df \right)^{-1},$  (17)

- Peak wavenumber  $k_p$ , computed via the peak period as in ref. 54:

$$\bar{T}_p = \frac{\int S(f)^4 \, df}{\int f \cdot S(f)^4 \, df}. \tag{18}$$

**Table 2. Hyperparameters used in experiments**

Hyperparameters	
Optimizer	Adam
Learning rate	$1 \times 10^{-4}$
Number of hidden layers	3
Neurons in hidden layers	$(32/\sqrt{n_h}, 16/\sqrt{n_h}, 8/\sqrt{n_h})$
$\ell_2$ penalty $\lambda_2$	$1 \times 10^{-5}$
Number of training epochs	50
Number of SWAG epochs	50
Number of SWAG posterior samples	100
Train-validation split	60% train, 40% test

$n_h$ : number of input heads.

Table 3. Full list of experiments

ID	Feature groups			Scores		
	1	2	3	$\mathcal{L} \times 10^4$	$\mathcal{E} \times 10^2$	$\mathcal{C} \times 10^2$
1	{ <i>r</i> }			4.51	8.23	3.35
2	{ <i>r</i> , <i>R</i> }	{Ur}		5.42	9.94	5.54
3	{ <i>r</i> , <i>R</i> , BFI}			5.43	10.50	5.60
4	{ <i>r</i> , <i>R</i> }	{Ur, <i>R</i> }		5.46	9.99	4.57
5	{ <i>r</i> , <i>R</i> }	{ $\epsilon$ , $\tilde{D}$ }		5.53	11.20	5.79
6	{ <i>r</i> , $\epsilon$ , $\tilde{D}$ }			5.20	11.00	3.60
7	{ <i>r</i> , $\epsilon$ , <i>R</i> }			5.31	11.40	6.97
8	{ <i>r</i> , $\tilde{D}$ , <i>R</i> }			5.41	11.50	7.31
9	{ $\epsilon$ , $\tilde{D}$ , <i>R</i> }			-0.13	24.80	7.60
10	{ $\sigma_f$ }	{ $\epsilon$ , $\tilde{D}$ , <i>R</i> }		3.93	13.60	9.02
11	{ <i>r</i> }	{ $\epsilon$ , $\tilde{D}$ , <i>R</i> }		5.41	10.60	7.18
12	{ <i>r</i> }	{ $\epsilon$ , $\tilde{D}$ }	{BFI, <i>R</i> }	5.41	11.10	6.02
13	{ <i>r</i> , <i>R</i> }	{ $\tilde{D}$ , $\epsilon$ , $\sigma_\theta$ }		5.99	12.40	4.06
14	{ <i>r</i> , <i>R</i> }	{ $\tilde{D}$ , $\epsilon$ , $\sigma_f$ }		5.82	11.80	6.37
15	{ <i>r</i> , <i>R</i> }	{ $\tilde{D}$ , $\epsilon$ , <i>R</i> }		5.62	11.00	5.45
16	{ <i>r</i> , $\epsilon$ , $\tilde{D}$ , $\sigma_\theta$ }			5.83	11.60	5.94
17	{ <i>r</i> }	{ $\epsilon$ , $\tilde{D}$ }	{BFI, $\sigma_f$ , $\sigma_\theta$ }	5.85	11.80	6.40
18	{ <i>r</i> }	{ $\epsilon$ , $\tilde{D}$ , $\sigma_f$ , $\sigma_\theta$ }		6.06	11.30	4.43
19	{ <i>r</i> , $\epsilon$ , $\tilde{D}$ , <i>R</i> , $\lambda_p$ }			5.86	13.50	7.13
20	{ <i>r</i> , $\epsilon$ , $\tilde{D}$ , $\sigma_\theta$ , $\nu$ }			6.18	12.80	6.78
21	{ <i>r</i> , $\epsilon$ , $\tilde{D}$ , $\sigma_\theta$ , $\nu$ , $E_h$ }			6.19	14.10	6.71
22	{ <i>r</i> , $\epsilon$ , $\tilde{D}$ , $\sigma_\theta$ , $\sigma_f$ , $\nu$ , $E_h$ }			6.42	15.70	4.97
23	{ <i>r</i> , $\epsilon$ , $\tilde{D}$ , $\sigma_\theta$ , $\sigma_f$ , $E_h$ , BFI, <i>R</i> }			6.60	17.00	5.75
24	{ <i>r</i> , $\epsilon$ , $\tilde{D}$ , $\sigma_\theta$ , $\sigma_f$ , $E_h$ , $H_s$ , $\bar{T}$ , $\kappa$ , $\mu$ , $\lambda_p$ }			6.51	19.50	4.95

$\mathcal{L}$ , Prediction score (higher is better);  $\mathcal{E}$ , Invariance error (lower is better);  $\mathcal{C}$ , Calibration error (lower is better); Color coding ranges between (median – IQR, median + IQR) with interquartile range IQR, with dark green indicating best and dark red worst. Symbols: *r*, Crest–trough correlation;  $\nu$ , Spectral bandwidth (narrowness);  $\sigma_f$ , Spectral bandwidth (peakedness);  $\sigma_\theta$ , Directional spread;  $\epsilon$ , Peak steepness  $H_s k_p$ ; *R*, Directionality index  $\sigma_\theta^2/(2\nu^2)$ ; BFI, Benjamin–Feir index;  $\tilde{D}$ , Relative peak water depth  $Dk_p/(2\pi)$ ;  $E_h$ , Relative high-frequency energy; Ur, Ursell number;  $\bar{T}$ , Mean period;  $\kappa$ , Kurtosis;  $\mu$ , Skewness;  $H_s$ , Significant wave height.

This leads to the peak wavenumber through the dispersion relation for linear waves in intermediate water of depth *D*:

$$f(k)^2 = \frac{gk}{(2\pi)^2} \tanh(kD). \tag{19}$$

- An approximate inverse is given in Fenton (55).
- Relative depth, based on the wavelength  $\lambda$ :

$$\tilde{D} = \frac{D}{\lambda} = \frac{1}{2\pi} k_p D, \tag{20}$$

- Peak steepness:

$$\epsilon = H_s k_p, \tag{21}$$

- Benjamin–Feir index:

$$\text{BFI} = \frac{\epsilon \nu}{\sigma_f} \sqrt{\max\{\beta/\alpha, 0\}}, \tag{22}$$

where  $\nu$ ,  $\alpha$ ,  $\beta$  are coefficients depending only on  $\tilde{D}$  full expression given in ref. 56.

- Directionality index:

$$R = \frac{\sigma_\theta^2}{2\nu_f^2}, \tag{23}$$

- Crest–trough correlation:

$$r = \frac{1}{m_0} \sqrt{\rho^2 + \lambda^2}, \tag{24}$$

$$\rho = \int_0^\infty \mathcal{S}(\omega) \cos\left(\omega \frac{\bar{T}}{2}\right) d\omega, \tag{25}$$

$$\lambda = \int_0^\infty \mathcal{S}(\omega) \sin\left(\omega \frac{\bar{T}}{2}\right) d\omega, \tag{26}$$

where  $\omega$  is the angular frequency and  $\bar{T} = m_0/m_1$  the spectral mean period (12).

**Model Implementation and Hyperparameters.** All performance critical model code is implemented in JAX (57), using neural network modules from flax (58) and optimizers from optax (59). We run each experiment on a single Tesla P100 GPU in about 40 min, including SWAG sampling and retraining on every validation subset. The whole training process can also be executed on CPU in about 2 h. The hyperparameters for all experiments are shown in Table 2.

**Full List of Experiments.** See Table 3.

**Reference Wave Height Distributions.** We use the following theoretical wave height exceedance distributions for comparison (with rogue wave threshold  $\kappa$ , here  $\kappa = 2$ ):

• Rayleigh (38): 
$$P_R(\kappa) = \exp(-2\kappa^2), \tag{27}$$

• Tayfun (12, 22): 
$$P_T(\kappa) = \exp\left(\frac{-4}{1+r}\kappa^2\right), \tag{28}$$

- Mori & Janssen (50, 60):

$$P_{MJ}(\kappa) = \left(1 + \frac{2\pi}{3\sqrt{3}} \frac{\text{BFI}^2}{1 + 7.1R} \kappa^2 (\kappa^2 - 1)\right) \exp(-2\kappa^2), \quad [29]$$

- Hybrid:

$$P_H(\kappa) = \left(1 + \frac{2\pi}{3\sqrt{3}} \frac{\text{BFI}^2}{1 + 7.1R} \kappa^2 (\kappa^2 - 1)\right) \exp\left(\frac{-4}{1+r} \kappa^2\right). \quad [30]$$

**Data, Materials, and Software Availability.** The preprocessed and aggregated version of the Free Ocean Wave Dataset, Coastal Data Information Program data used in this study is available for download at <https://erda.ku.dk/archives/ee6b452c1907fbd48271b071c3cee10e/published-archive.html> (61).

All model code is openly available at <https://github.com/dionhaefner/rogue-wave-discovery> (62). This publication was made possible by the following opensource software stack: JAX (57), flax (58), optax (59), PySR (9), scikit-learn (63), PyALE (64), NumPy (65), SciPy (66), matplotlib (67), Seaborn (68), pandas (69), and Jupyter (70).

**ACKNOWLEDGMENTS.** D.H. received funding from the Danish Offshore Technology Centre. Raw data were furnished by the Coastal Data Information Program, Integrative Oceanography Division, operated by the Scripps Institution of Oceanography, under the sponsorship of the United States Army Corps of Engineers and the California Department of Parks and Recreation. Computational resources were provided by DC<sup>3</sup>, the Danish Center for Climate Computing. Portions of this work were developed from the doctoral thesis of D.H. (71). We thank Jonas Peters for helpful discussions in the early stages of this work. We thank two anonymous reviewers for their helpful comments.

- E. Didenkulova, Catalogue of rogue waves occurred in the World Ocean from 2011 to 2018 reported by mass media sources. *Ocean Coast. Manag.* **188**, 105076 (2019).
- J. M. Dudley, G. Genty, A. Mussot, A. Chabchoub, F. Dias, Rogue waves and analogies in optics and oceanography. *Nat. Rev. Phys.* **1**, 675–689 (2019).
- J. Behrens, J. Thomas, E. Terrill, R. Jensen, "CDIP: Maintaining a robust and reliable ocean observing buoy network" in 2019 IEEE/OES Twelfth Current, Waves and Turbulence Measurement (CWTM) (2019), pp. 1–5.
- D. Häfner, J. Gemmrich, M. Jochum, Real-world rogue wave probabilities. *Sci. Rep.* **11**, 10084 (2021).
- A. Cattrell, M. Srokosz, B. Moat, R. Marsh, Can rogue waves be predicted using characteristic wave parameters? *J. Geophys. Res. Oceans* **123**, 5624–5636 (2018).
- D. Häfner, J. Gemmrich, M. Jochum, FOWD: A free ocean wave dataset for data mining and machine learning. *J. Atmos. Oceanic Technol.* **1**, 1305–1322 (2021).
- C. Molnar, *Interpretable Machine Learning: A Guide for Making Black Box Models Explainable, Second Edition* (2020). <https://christophm.github.io/interpretable-ml-book/cite.html>.
- M. Cranmer *et al.*, Discovering symbolic models from deep learning with inductive biases. *NeurIPS* **2020**, 17429–17442 (2020).
- M. Cranmer, Interpretable machine learning for science with PySR and SymbolicRegression.jl. *arXiv [Preprint]* (2023). <https://arxiv.org/abs/2305.01582> (Accessed 25 September 2023).
- E. O. Voit, Perspective: Dimensions of the scientific method. *PLoS Comput. Biol.* **15**, e1007279 (2019).
- T. A. A. Adcock, P. H. Taylor, The physics of anomalous ("rogue") ocean waves. *Rep. Progr. Phys.* **77**, 105901 (2014).
- M. A. Tayfun, F. Fedele, Wave-height distributions and nonlinear effects. *Ocean Eng.* **34**, 1631–1649 (2007).
- M. Miche, Mouvements ondulatoires de la mer en profondeur constante ou décroissante. *Annales de Ponts et Chaussées*, 1944, pp(1) 26–78, (2) 270–292, (3) 369–406 (1944).
- J. Gemmrich, C. Garrett, Dynamical and statistical explanations of observed occurrence rates of rogue waves. *Nat. Hazards Earth Syst. Sci.* **11**, 1437–1446 (2011).
- F. Fedele, J. Brennan, S. Ponce de León, J. Dudley, F. Dias, Real world ocean rogue waves explained without the modulational instability. *Sci. Rep.* **6**, 27715 (2016).
- M. Onorato *et al.*, Extreme waves, modulational instability and second order theory: Wave flume experiments on irregular waves. *Eur. J. Mech. - B/Fluids* **25**, 586–601 (2006).
- K. Trulsen, H. Zeng, O. Gramstad, Laboratory evidence of freak waves provoked by non-uniform bathymetry. *Phys. Fluid.* **24**, 097101 (2012).
- J. K. Mallory, Abnormal waves on the south east coast of South Africa. *Int. Hydrogr. Rev.* **51** (1974).
- E. G. Didenkulova, T. G. Talipova, E. N. Pelinovsky, "Rogue waves in the drake passage: Unpredictable hazard" in *Antarctic Peninsula Region of the Southern Ocean: Oceanography and Ecology, Advances in Polar Ecology*, E. G. Morozov, M. V. Flint, V. A. Spiridonov, Eds. (Springer International Publishing, Cham, 2021), pp. 101–114.
- M. L. McAllister, S. Draycott, T. A. A. Adcock, P. H. Taylor, T. S. van den Bremer, Laboratory recreation of the Draupner wave and the role of breaking in crossing seas. *J. Fluid Mech.* **860**, 767–786 (2019).
- J. Pearl, *Causality* (Cambridge University Press, 2009).
- M. A. Tayfun, Distribution of large wave heights. *J. Water. Port Coast Ocean Eng.* **116**, 686–707 (1990).
- Y. Goda, Reanalysis of regular and random breaking wave statistics. *Coast. Eng. J.* **52**, 71–106 (2010).
- T. Tang, D. Barratt, H. B. Bingham, T. S. van den Bremer, T. A. Adcock, The impact of removing the high-frequency spectral tail on rogue wave statistics. *J. Fluid Mech.* **953**, A9 (2022).
- D. J. Korteweg, G. De Vries, On the change of form of long waves advancing in a rectangular canal, and on a new type of long stationary waves. *London Edinbur. Dublin Philos. Magaz. J. Sci.* **39**, 422–443 (1895).
- P. Janssen, Shallow-water version of the Freak Wave Warning System, (ECMWF). Technical memorandum 813 (2018).
- P. A. E. M. Janssen, Nonlinear four-wave interactions and freak waves. *J. Phys. Oceanogr.* **33**, 863–884 (2003).
- F. Ursell, The long-wave paradox in the theory of gravity waves. *Math. Proc. Camb. Philos. Soc.* **49**, 685–694 (1953).
- ECMWF, Part VII, ECMWF Wave model in IFS Documentation CY47R3, IFS Documentation (ECMWF, 2021).
- K. Hornik, Approximation capabilities of multilayer feedforward networks. *Neural Netw.* **4**, 251–257 (1991).
- V. Nair, G. E. Hinton, "Rectified linear units improve restricted Boltzmann machines" in *Proceedings of the 27th International Conference on Machine Learning, ICML 2010* (OmniPress, Madison, WI, USA, 2010), pp. 807–814.
- D. P. Kingma, J. Ba, Adam: A method for stochastic optimization. *arXiv [Preprint]* (2014). <http://arxiv.org/abs/1412.6980> (Accessed 25 September 2023).
- W. Maddox, T. Garipov, P. Izmailov, D. Vetrov, A. G. Wilson, A simple baseline for Bayesian uncertainty in deep learning. *arXiv [Preprint]* (2019). <http://arxiv.org/abs/1902.02476> (Accessed 25 September 2023).
- J. Peters, P. Bühlmann, N. Meinshausen, Causal inference by using invariant prediction: Identification and confidence intervals. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **78**, 947–1012 (2016).
- J. Peters, D. Janzing, B. Schölkopf, "Elements of causal inference: Foundations and learning algorithms" in *Adaptive Computation and Machine Learning Series*, F. Bach, Ed. (MIT Press, Cambridge, MA, USA, 2017).
- C. Heinze-Deml, J. Peters, N. Meinshausen, Invariant causal prediction for nonlinear models. *J. Causal Inf.* **6**, 20170016 (2018).
- P. Xenopoulos, J. Rulff, L. G. Nonato, B. Barr, C. Silva, Calibrate: Interactive analysis of probabilistic model output. *IEEE Trans. Visual. Comput. Graphics* **29**, 853–863 (2022).
- M. S. Longuet-Higgins, On the statistical distribution of the height of sea waves. *JMR* **11**, 245–266 (1952).
- M. Schmidt, H. Lipson, Distilling free-form natural laws from experimental data. *Science* **324**, 81–85 (2009).
- S. M. Udrescu, *et al.*, Al Feynman 2.0: Pareto-optimal symbolic regression exploiting graph modularity. *Adv. Neural Inf. Process. Syst.* **33**, 4860–4871 (2020).
- S. L. Brunton, J. L. Proctor, J. N. Kutz, Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proc. Natl. Acad. Sci. U.S.A.* **113**, 3932–3937 (2016).
- K. R. Broløs *et al.*, An approach to symbolic regression using Feyn. *arXiv [Preprint]* (2021). <http://arxiv.org/abs/2104.05417> (Accessed 25 September 2023).
- J. H. Holland, Genetic algorithms. *Sci. Am.* **267**, 66–73 (1992).
- F. Fedele, M. A. Tayfun, On nonlinear wave groups and crest statistics. *J. Fluid Mech.* **620**, 221–239 (2009).
- F. Fedele, J. Herterich, A. Tayfun, F. Dias, Large nearshore storm waves off the Irish coast. *Sci. Rep.* **9**, 15406 (2019).
- O. Gramstad, K. Trulsen, Influence of crest and group length on the occurrence of freak waves. *J. Fluid Mech.* **582**, 463–472 (2007).
- W. Xiao, Y. Liu, G. Wu, D. K. P. Yue, Rogue wave occurrence and dynamics by direct simulations of nonlinear wave-field evolution. *J. Fluid Mech.* **720**, 357–392 (2013).
- J. Gemmrich, J. Thomson, Observations of the shape and group dynamics of rogue waves. *Geophys. Res. Lett.* **44**, 1823–1830 (2017).
- D. W. Apley, J. Zhu, Visualizing the effects of predictor variables in black box supervised learning models. *arXiv [Preprint]* (2019). <http://arxiv.org/abs/1612.08468> (Accessed 25 September 2023).
- N. Mori, M. Onorato, P. A. Janssen, On the estimation of the kurtosis in directional sea states for freak wave forecasting. *J. Phys. Oceanogr.* **41**, 1484–1497 (2011).
- L. H. Ying, Z. Zhuang, E. J. Heller, L. Kaplan, Linear and nonlinear rogue wave statistics in the presence of random currents. *Nonlinearity* **24**, R67–R87 (2011).
- M. Casas-Prat, L. H. Holthuijsen, Short-term statistics of waves observed in deep water. *J. Geophys. Res.: Oceans* **115** (2010).
- J. Gemmrich, L. Cicon, Generation mechanism and prediction of an observed extreme rogue wave. *Sci. Rep.* **12**, 1–10 (2022).
- I. R. Young, The determination of confidence limits associated with estimates of the spectral peak frequency. *Ocean Eng.* **22**, 669–686 (1995).
- J. D. Fenton, The numerical solution of steady water wave problems. *Comput. Geosci.* **14**, 357–368 (1988).
- M. Serio, M. Onorato, A. Ra. Osborne, P. Janssen, On the computation of the Benjamin-Feir index. *Nuovo Ciment. della Soc. Ital. Fisica C* **28**, 893–903 (2005).
- J. Bradbury *et al.*, JAX: composable transformations of Python+NumPy programs (Version 0.2.5, 2018). <http://github.com/google/jax>. Accessed 25 September 2023.
- J. Heek *et al.*, Flax: A neural network library and ecosystem for JAX (Version 0.4.0, 2020). <http://github.com/google/flax>. Accessed 25 September 2023.
- M. Hessel *et al.*, Optax: Composable gradient transformation and optimisation, in JAX! (Version 0.0.1, 2020). <http://github.com/deepmind/optax>. Accessed 25 September 2023.

60. N. Mori, P. A. E. M. Janssen, On kurtosis and occurrence probability of freak waves. *J. Phys. Oceanogr.* **36**, 1471–1483 (2006).
61. D. Häfner, Big Data Big Waves - Data files. Electronic Research Data Archive (ERDA), University of Copenhagen. <https://erda.ku.dk/archives/ee6b452c1907fbd48271b071c3cee10e/published-archive.html>. Deposited 14 August 2023.
62. D. Häfner, dionhaefner/rogue-wave-discovery: Code for the paper "Machine-Guided Discovery of a Real-World Rogue Wave Model" (2023). GitHub. <https://github.com/dionhaefner/rogue-wave-discovery>. Deposited 20 September 2023.
63. F. Pedregosa *et al.*, Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.* **12**, 2825–2830 (2011).
64. D. Jomar, PyALE: A Python implementation of accumulated local effect plots (Version 1.1.2, 2020). <https://github.com/DanaJomar/PyALE>.
65. C. R. Harris *et al.*, Array programming with NumPy. *Nature* **585**, 357–362 (2020).
66. P. Virtanen *et al.*, SciPy 1.0: Fundamental algorithms for scientific computing in Python. *Nat. Methods* **17**, 261–272 (2020).
67. J. D. Hunter, Matplotlib: A 2D graphics environment. *Comput. Sci. Eng.* **9**, 90–95 (2007).
68. M. L. Waskom, Seaborn: Statistical data visualization. *J. Open Source Softw.* **6**, 3021 (2021).
69. W. McKinney, "Data structures for statistical computing in Python" in *Proceedings of the 9th Python in Science Conference*, S. van der Walt, J. Millman, Eds. (2010), pp. 56–61.
70. T. Kluyver *et al.*, "Jupyter notebooks - a publishing format for reproducible computational workflows" in *Positioning and Power in Academic Publishing: Players, Agents and Agendas*, F. Loizides, B. Schmidt, Eds. (IOS Press, Netherlands, 2016), pp. 87–90.
71. D. Häfner, "An Ocean of Data: Inferring the Causes of Real-World Rogue Waves," PhD thesis, Niels Bohr Institute, Faculty of Science, University of Copenhagen (2022).