

Can Machine Learning Predict Extreme Events in Complex Systems?

Di Qi^{a,1,2} and Andrew J. Majda^{a,1,2}

^aDepartment of Mathematics and Center for Atmosphere and Ocean Science, Courant Institute of Mathematical Sciences, New York University, New York, NY 10012

This manuscript was compiled on October 3, 2019

Extreme events and the related anomalous statistics are ubiquitously observed in many natural systems, and it remains a grand challenge for the development of efficient methods to understand and accurately predict such representative features. Here, we investigate the skill of deep learning strategies in the prediction of extreme events in complex turbulent dynamical systems. Deep neural networks have been successfully applied to many imaging processing problems involving big data, and have recently shown potential for the study of dynamical systems. We propose to use a densely connected mixed-scale network model to capture the extreme events appearing in a truncated KdV (tKdV) statistical framework which creates anomalous skewed distributions consistent with recent laboratory experiments for shallow water waves across an abrupt depth change, where a remarkable statistical phase transition is generated by varying the inverse temperature parameter in the corresponding Gibbs invariant measures. The neural network is trained using data without knowing the explicit model dynamics, and the training data is only drawn from the near-Gaussian regime of the tKdV model solutions without the occurrence of large extreme values. A relative entropy loss function together with empirical partition functions is proposed for measuring the accuracy of the network output where the dominant structures in the turbulent field are emphasized. The optimized network is shown to gain uniformly high skill in accurately predicting the solutions in a wide variety of statistical regimes including highly skewed extreme events. The new technique is promising to be further applied to other complicated high-dimensional systems.

anomalous extreme events | convolutional neural networks | turbulent dynamical systems

Extrême events and their anomalous statistics are ubiquitous in various complex turbulent systems such as the climate, material and neuroscience, as well as engineering design (1–4). Understanding and accurate prediction of such phenomena remain a grand challenge, and have become an active contemporary topic in applied mathematics (5–8). Extreme events can be isolated rare events (2, 9, 10), or they can often be intermittent and even frequent in space and time (6, 8, 11, 12). The curse of dimension forms one important obstacle for the accurate prediction of extreme events in large complex systems (3, 4, 6, 13), where both novel models and efficient numerical algorithms are required. A typical example can be found in recent laboratory experiments for turbulent surface water waves going through an abrupt depth change revealing a remarkable transition to anomalous extreme events from near-Gaussian incoming flows (1).

A statistical dynamical model is then proposed in (14, 15) that successfully predicts the anomalous extreme behaviors observed in the shallow water wave experiments. The truncated Korteweg-de Vries (tKdV) equation is proposed as the governing equation modeling the flow surface displacement. Gibbs invariant measures are induced based on the Hamilto-

nian form of the tKdV equation to describe the probability distributions at equilibrium. A statistical transition from symmetric near-Gaussian statistics to a highly skewed probability density function (PDF) is achieved by simply controlling the ‘inverse temperature’ parameter in the Gibbs measure (15).

In recent years, machine learning strategies, particularly the deep neural networks, have been extensively applied to a wide variety of problems involving big data, such as image classification and identification (16–19). On the other hand, it still remains an actively growing topic to construct proper deep learning strategies for the study of complex turbulent flows. The deep neural network tools developed for imaging processing have been suggested to be applied for data-driven predictions of chaotic dynamical systems (20, 21), climate and weather forecasts (22, 23), and parameterization of unresolved processes (24–26). In the statistical prediction of extreme events, the available data for training is often restricted in limited regimes. A successful neural network is required to maintain adaptive skill in wider statistical regimes with vastly distinct statistics away from the training dataset. Besides, a working prediction model for turbulent systems would also require the prediction time scale longer than the decorrelation time that characterizes the mixing rate of the state variables.

In this paper, we investigate the extent of skill of the deep neural networks in predicting statistical solutions of complex turbulent systems, especially involving highly skewed probability density functions. The statistical tKdV equations serve as a difficult first test model for extreme event prediction with

Significance Statement

Understanding and predicting extreme events as well as the related anomalous statistics is a grand challenge in complex natural systems. Deep convolutional neural networks provide a useful tool to learn the essential model dynamics directly from data. A new deep learning strategy is proposed to predict the extreme events appeared in turbulent dynamical systems. A truncated KdV model displaying distinctive statistics from near-Gaussian to highly skewed distributions is used as the test model. The neural network is trained using data only from the near-Gaussian regime without the occurrence of large extreme values. The optimized network demonstrates uniformly high skill in successfully capturing the solution structures in a wide variety of statistical regimes including the highly skewed extreme events.

D. Q. and A. J. M. designed the research, performed the research, and wrote the paper.

The authors declare no conflict of interest.

¹D. Q. and A. J. M. contributed equally to this work.

²To whom correspondence should be addressed. E-mail: qidi@cims.nyu.edu, jonjon@cims.nyu.edu.

simple trackable dynamics but a rich variety of statistical regimes from near Gaussian to highly skewed PDFs showing extreme events. The important questions to ask are whether the deep networks can be trained to learn the complex hidden structures in the highly nonlinear dynamics purely from restricted data, and what are the essential structures required in the network to gain the ability to capture extreme events.

Our major goal here is to get accurate statistical prediction for the extreme events in time intervals significantly longer than the decorrelation time of the complex turbulent system. To achieve this, a convolutional neural network (MS-D Net) which exploits multi-scale connections and densely connected structures (27) is adopted to provide the basic network architecture to be trained using the model data from the tKdV equation solutions. This network enjoys the benefits of simpler model implementation and a smaller number of tunable training hyperparameters. Thus it becomes much easier to train requiring less computational cost and technical tuning of the hyperparameters.

The key structures for the neural network to successfully capture extreme events include: i) the use of a relative entropy (*Kullback-Leibler* divergence) loss function to calibrate the closeness between the target and the network output as distribution functions, so that the crucial shape of the model solutions are captured instead of a pointwise fitting in the turbulent output field values; and ii) calibrating the output data under a combination of empirical partition functions emphasizing the large positive and negative values in the model prediction so that the main features in the solutions are further emphasized. This convolutional neural network model enjoys the following major advantages for the prediction of extreme events:

- The simple basic network architecture makes the model easier to train and efficient to predict the solutions among different statistical scenarios, as well as reduces the danger of overfitting from data;
- The network structure approximates the original model dynamics with the designed model loss function and treatment of output data. So the complex system dynamics is easier to be learned from data;
- The temporal and spatial correlations in different scales are modeled automatically from the design of the network with convolution kernels representing different scales in different layers;
- The method shows robust performance with different model hyperparameters, and can be generalized for the prediction of more complicated turbulent systems.

Direct numerical tests show high skill of the neural network in successfully capturing the extreme values in the solutions with the model parameters only learned from the near-Gaussian regime of vastly different statistics. The model also displays accurate prediction in much longer time beyond the decorrelation time scale of the state, proving the robustness of the methods. The successful prediction in the tKdV equation implies the potential of future applications of the network to more complicated high dimensional systems.

Background for extreme events and neural network structures

The truncated KdV equations with extreme events. The tKdV model provides a desirable set of equations capable of capturing many complex features in surface water wave turbulence with simple trackable dynamics. Through a high wavenumber cutoff at Λ (with $J = 2\Lambda + 1$ grid points), the Galerkin projected state $u = \sum_{1 \leq |k| \leq \Lambda} \hat{u}_k(t) e^{ikx}$ induces stronger turbulent dynamics than the original continuous one (15). The tKdV equation has been adapted to describe the sudden phase transition in statistics (14) where highly skewed extreme events are generated from near-Gaussian statistics for waves propagating across an abrupt depth change. The tKdV model is formulated on a periodic domain $x \in [-\pi, \pi]$ as

$$u_t + E_0^{1/2} L_0^{-3/2} D_0^{-3/2} u u_x + L_0^{-3} D_0^{1/2} u_{xxx} = 0, \quad [1]$$

where the state variable $u(x, t)$ represents the surface wave displacement to be learned directly using the deep neural network. The model is nondimensionalized using the characteristic scales E_0 as the total energy, L_0 as the length scale, and D_0 as the water depth. The steady state distribution of the tKdV solution can be described by the invariant Gibbs measure derived from the equilibrium statistical mechanics (28)

$$\mathcal{G}_\theta(u) = C_\theta^{-1} \exp \left\{ -\theta \left[E_0^{1/2} L_0^{-3/2} D_0^{-3/2} H_3(u) - L_0^{-3} D_0^{1/2} H_2(u) \right] \right\} \delta(E(u) - 1), \quad [2]$$

with a competition between the cubic term $H_3 = \frac{1}{6} \int u^3$ and the quadratic term $H_2 = \frac{1}{2} \int u_x^2$ from the Hamiltonian. The last term in [2] is the delta function constraining the total energy conservation, $E(u) = \frac{1}{2} \int u^2 = 1$. The only parameter $\theta < 0$ as the ‘inverse temperature’ determines the skewness of the PDF of u (14). The Gibbs measures [2] with different values of θ can be used to provide initial samples for the direct simulations of the model [1]. Different final equilibrium statistics (with various skewness) can be obtained based on the initial configuration of the ensemble set (that is, from picking different inverse temperature values θ). A detailed description about the statistical tKdV model together with the simulation setup is provided in *SI Appendix, A*.

Training and prediction data from the same model with distinct statistics. The basic idea in training the deep neural network is to use a training dataset with solutions sampled from [2] using near-Gaussian statistics. The tKdV model dynamics can be learned from the training set without explicitly knowing the model dynamics. Then the question is what is the range of skill in the trained neural network to predict the highly skewed non-Gaussian distributions among different sets of data. The training and prediction datasets are proposed from the ensemble solutions of the tKdV model [1] based on the following strategy:

- In the training dataset, we generate solutions from an ensemble simulation starting from a near-Gaussian PDF using a small absolute value of the inverse temperature θ_0 (as shown in the first row and the near-Gaussian PDF in Figure 1). On one hand, the model dynamics is represented by the group of solutions $\{u_{\theta_0}\}$ to be learned

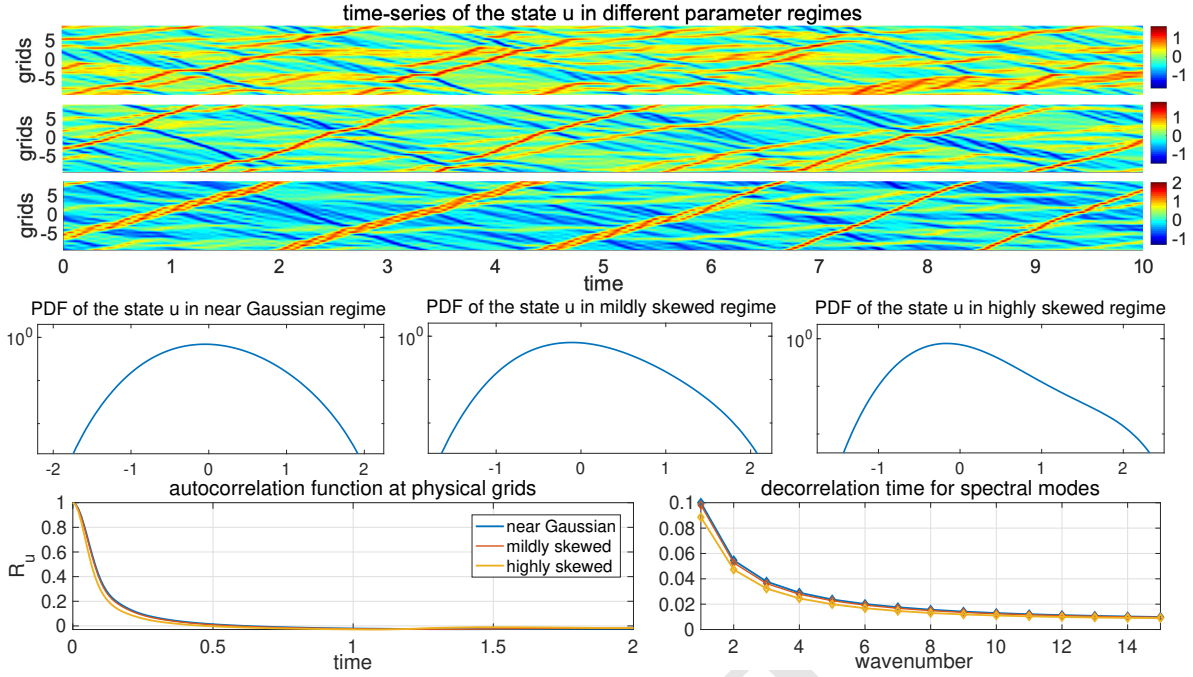


Fig. 1. Solutions and statistics of the tKdV equation in three typical parameter regimes with different statistics. The first three rows plot solution trajectories in the three regimes with near-Gaussian (first row), mildly skewed (second row), and highly skewed (third row) statistics. The corresponding equilibrium PDFs of the three cases are shown next. The autocorrelation functions and decorrelation time of each Fourier mode of the model state u are compared in the last row.

161 through the deep neural network. On the other hand, 162 only near-Gaussian statistics is obtained in this training 163 dataset so the neural network cannot know about the 164 skewed rare events appearing in other statistical regimes 165 directly from the training process.

- 166 • For the prediction dataset, we test the model skill using 167 the data $\{u_\theta\}$ generated from various different initial 168 inverse temperatures θ (as shown in the second and third 169 rows as well as the skewed PDFs in Figure 1). It provides 170 an interesting testbed to check the scope of skill in the 171 optimized neural network for capturing the distinctive 172 statistics and extreme events.

173 The choices of the training and prediction datasets are illus- 174 trated in Figure 1, which first shows in the first three rows 175 realizations of the tKdV model solutions from different inverse 176 temperatures θ . A smaller amplitude of θ gives near-Gaussian 177 statistics in the model state u , while larger amplitudes of 178 θ give trajectories with skewed PDFs. The corresponding 179 equilibrium PDFs of the state u from direct ensemble dynam- 180 ical solutions are compared next to illustrate explicitly the 181 transition in statistics. Turbulent dynamics with multiscale 182 structures are observed in all the tKdV solutions. The autocor- 183 relation function $\mathcal{R}_u(t) = \langle u(t)u(0) \rangle$ and the decorrelation 184 time $T_{\text{decorr}} = \int \mathcal{R}_u(t)$ are plotted in the bottom row of Figure 185 1, confirming the rapid mixing in the solutions.

186 **Data structures for the deep neural networks.** The deep convolu- 187 tional networks can be viewed as a function $\mathbf{y} = f_M(\mathbf{x})$, map- 188 ping the input signal $\mathbf{x} \in \mathbb{R}^{m \times n \times c}$ with m rows, n columns, and 189 c channels to the output data $\mathbf{y} \in \mathbb{R}^{m' \times n' \times c'}$ with m' rows, n' 190 columns, and c' channels. We consider an ensemble simulation 191 of M trajectories of the tKdV equation evaluated at the J grid

192 points and measured in a time interval $[t_0, t_{N-1}]$. Thus the input 193 data for the network comes from the ensemble solutions at 194 the first N time measurements at $t_j = j\Delta t, j = 0, \dots, N-1$,

$$\mathbf{x}^{(l)} = [u_0^{(l)}, u_1^{(l)}, \dots, u_{J-1}^{(l)}]^T (t_0, \dots, t_{N-1}) \in \mathbb{R}^{J \times N}, \quad 195$$

196 as a tensor with $m = J$ rows for the spatial grid points, $n = N$ 197 columns for the discretized time evaluations, and only one 198 channel $c = 1$ for each of the input samples $l = 1, \dots, M$. An 199 ensemble of total M independent solutions from the Monte- 200 Carlo simulation is divided into mini-batches to feed in the 201 network in the training process. For simplicity, the output 202 data is designed in the same shape as the predicted states 203 evaluated at a later time $t = T + t_0$ starting from the previous 204 initial data

$$\mathbf{y}^{(l)} = [u_0^{(l)}, u_1^{(l)}, \dots, u_{J-1}^{(l)}]^T (T + t_0, \dots, T + t_{N-1}) \in \mathbb{R}^{J \times N}. \quad 205$$

206 The forwarding time T controls how long we would like the 207 network to push forward the states u in one time update. 208 For one effective neural network for the complex system, the 209 time scale T is expected to be longer than the decorrelation 210 time $T > T_{\text{decorr}}$. The above construction is supposed to 211 feed both the time and spatial correlations of the original 212 dynamical model into the neural network to be learned in the 213 approximation map $\mathbf{y} = f_M(\mathbf{x})$.

214 **Deep convolutional neural network architecture.** The basic 215 structures of the convolutional neural network include the 216 operators in each single convolution layer; and the connections 217 between multiple layers. We would like to first keep the neural 218 network in its simplest standard setup, so that we are able 219 to concentrate on the improvement in key structures without 220 risking at getting lost in manipulating the various complicated

221 *ad hoc* hyperparameters. More detailed convolutional network
 222 construction is described in *SI Appendix, B.1* following the
 223 general neural network architecture as in (19, 27).

224 **Basic convolutional neural network unit.** In each single convolution
 225 layer, the input data from the previous layer output is updated
 226 in the general form

$$227 \quad \mathbf{y} = \sigma(g_h(\mathbf{x}) + b).$$

228 A convolutional operator g_h is first applied on the input data
 229 \mathbf{x} in a small symmetric window with size $w \times w$, where the
 230 first dimension controls the correlation in the spatial direction
 231 and the second one for the temporal correlation. A bias b is
 232 added to the convolved data before applying a final nonlinear
 233 operator σ using the common choice of rectified linear unit
 234 (ReLU) function. The convolution kernel starts with a small
 235 size 3×3 (that is, using only the two nearest neighbor points
 236 in space and time) which enables fast computation and easy
 237 control. Naturally, periodic boundary condition is applied
 238 on the spatial dimension and replicate boundary is added in
 239 time before t_0 and after t_{N-1} for the boundary padding. No
 240 additional structures are implemented in the convolution layer
 241 unit to keep the basic standard architecture used in imaging
 242 processing (19).

243 **A densely connected and mix-scale structure.** Next, we need to
 244 propose the connection between different layers. The common
 245 *feedforward deep neural network* feeds the input data in the
 246 i -th layer only to the next $(i + 1)$ -th layer. The feedforward
 247 network often requires a larger number of layers to work,
 248 thus it is expansive to train and difficult to handle. Proper
 249 downscaling and upscaling steps going through the layers
 250 may also be required, while these downscaling and upscaling
 251 operations may not be a feasible approach for simulating the
 252 dynamical model time integration steps.

253 As an alternative approach, a *mixed-scale dense neural*
 254 *network* (MS-D Net) is introduced in (27) by mixing differ-
 255 ent scales within each layer using a dilated convolution, and
 256 densely connecting all the feature maps among all the layers.
 257 First, to learn the multiscale structures, the convolution ker-
 258 nels in different layers are dilated differently by adding s zeros
 259 between the values in the original kernel $w \times w$. The dilated
 260 convolutions become especially appealing for capturing the
 261 multiscale structures in the turbulent dynamics. Different spa-
 262 tial and temporal scales are included adaptively with different
 263 convolution length scales. Second, the dense network connec-
 264 tion includes all the previous layer information to update the
 265 output data in the next layer. In the implementation, all the
 266 previous layer outputs are piled together as input channels
 267 for the next layer. Together with the multiscale convolution
 268 kernels used in different layers, the output in the next layer
 269 combines the information in different scales and produces a
 270 balanced update in the next step.

271 The mixed-scale dense neural network requires fewer fea-
 272 ture maps and trainable parameters, so it is easier to handle
 273 compared with the direct feedforward network. It provides
 274 a desirable setup for the prediction in dynamical systems by
 275 feeding in all the data in previous layers decomposed into
 276 different scale structures. Then information at different scales
 277 communicates with each other through the dense network
 278 connection. Intuitively, this is a reasonable structure for the
 279 turbulent solutions since all the history information is useful

280 for the prediction in the next steps. The densely connected
 281 network structure is also comparable to the time integration
 282 scheme, where all the history information is used to update
 283 the state at the next time step without using any downscaling
 284 and upscaling steps.

285 A new learning strategy for extreme event prediction

286 In this section, we construct the specific network structures
 287 designed for learning turbulent system dynamics then the
 288 prediction of extreme events. In this case with data from
 289 the turbulent models, small fluctuations in the solutions may
 290 introduce large errors in optimizing the loss function. We aim
 291 at capturing the dominant emerging features such as extreme
 292 events and are more interested in the statistical prediction
 293 rather than the exact locations of the extreme values.

294 **Loss function calibrating the main dynamical features.** In the
 295 training process, the model parameters are achieved through
 296 the optimization for the proper loss function (or cost) proposed
 297 depending on the target to be captured using the network.
 298 The two popular standard choices for the loss functions are
 299 using the L_1 error and the mean square L_2 error:

- L_1 error loss: this criterion measures the mean absolute
 error between each element in the output data \mathbf{x} and
 target \mathbf{y} through the L_1 distance

$$299 \quad l_1(\mathbf{x}, \mathbf{y}) = \frac{1}{M} \sum_{m=1}^M \|\mathbf{x}^{(m)} - \mathbf{y}^{(m)}\|, \quad [3] \quad 303$$

304 where M is the training data size in one cycle, and
 305 $(\mathbf{x}^{(m)}, \mathbf{y}^{(m)})$ is one member of output data and target
 306 data from the training mini-batch.

- L_2 error loss: this criterion measures the mean squared
 error between each element in the output \mathbf{x} and target \mathbf{y}
 through the mean square L_2 norm

$$307 \quad l_2(\mathbf{x}, \mathbf{y}) = \frac{1}{M} \sum_{m=1}^M \|\mathbf{x}^{(m)} - \mathbf{y}^{(m)}\|^2, \quad [4] \quad 310$$

311 among all the numbers of training samples M .

312 The above two loss functions offer pointwise measurements of
 313 the errors in space-time for each predicted sample from the
 314 network. This may cause problems especially when the system
 315 for prediction is highly turbulent with internal instability and
 316 is fast mixing. Small error perturbation in the input data
 317 may lead to vastly different solutions shortly after the mixing
 318 decorrelation time. The pointwise measurements focus on
 319 the accuracy at each value of the solutions, thus the small
 320 fluctuation errors might be accumulated and amplified under
 321 such metrics and unnecessary large weights are added to correct
 322 errors in the moderate-amplitude fluctuation parts.

323 On the other hand, we are most interested in the prediction
 324 of statistical features in the extreme events rather than the
 325 exact trajectory solutions of the system. The small shifts in the
 326 extreme value locations should be tolerated in the loss function.
 327 Therefore, a more useful choice could be the *relative entropy*
 328 loss function that measures the *Kullback-Leibler* divergence in
 329 the predicted density functions:

- *Relative entropy loss*: the relative entropy loss function computes the distance between two distribution functions

$$l_{\text{KL}}(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = \frac{1}{M} \sum_{m=1}^M \sum_i \tilde{y}_i^{(m)} \log \frac{\tilde{y}_i^{(m)}}{\tilde{x}_i^{(m)}}, \quad [5]$$

where the superscript m represents the mini-batch members to be measured in the relative entropy metric, and the subscript i goes through all the dimensions of the normalized variables $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ to be described next.

Under the relative entropy loss function in [5], the input data \mathbf{x} and \mathbf{y} are treated as distribution functions. The ‘shapes’ between the output data and the target are compared rather than the pointwise details, so that it guides the network to focus on the main model dynamical features instead of the turbulent fluctuations that are impossible to be fitted accurately. The additional difficulty to train the network using the relative entropy loss function is the constraint on the form of the output data to measure. The input of the relative entropy requires to be in the form of a density distribution function.

Scaling the output data with empirical partition functions. In this case of using the relative entropy loss function, we need to propose proper preprocessing of the output and target data to fit the required structure as a probability distribution. One direct way to do so is by taking the *softmax* function for the output data from the neural network normalized by a partition function

$$\tilde{x}_i = \frac{\exp(x_i)}{\sum_i \exp(x_i)}, \quad [6]$$

before measuring the error in the loss function. In this way, the data to put into the relative entropy loss function is normalized inside the range $[0, 1]$ with summation 1. This agrees with the definition in the relative entropy inputs. More importantly, this normalization emphasizes the large positive values of the data. Thus it offers a better calibration for the occurrence of positive major flow structures to be captured in the solutions.

Furthermore, a better choice for balancing both the positive and negative dominant values in the training data is to introduce scales with ‘temperatures’. We use the following two empirical partition functions with both positive and negative coefficients to rescale the output data as

$$\tilde{x}_i^+ = \frac{\exp(x_i/T_+)}{\sum_i \exp(x_i/T_+)}, \quad \tilde{x}_i^- = \frac{\exp(-x_i/T_-)}{\sum_i \exp(-x_i/T_-)}, \quad [7]$$

where $T_+ > 0, T_- > 0$ are the positive and negative temperatures weighing the importance of dominant large amplitude features in the scaled measure. Accordingly, the loss function to minimize under the relative entropy metric becomes a combination with the two empirical partition functions

$$l_{\text{EPF}}(\mathbf{x}, \mathbf{y}) = l_{\text{KL}}(\tilde{\mathbf{x}}^+, \tilde{\mathbf{y}}^+) + \alpha l_{\text{KL}}(\tilde{\mathbf{x}}^-, \tilde{\mathbf{y}}^-), \quad [8]$$

where we use $\alpha > 0$ as a further balance between the positive and negative temperature components. In this combined empirical partition function metric using [7] and [8], the major flow structures in the turbulent field represented by the dominant extreme values are better characterized from both the positive and negative sides in the statistics of the model. In the following computational experiments, we always pick $T_+ = T_- = 1$ and $\alpha = 1$ for simplicity. More discussion for role of the balance weight α is shown in *SI Appendix, B.2*.

Model performance in training and prediction stages

Using the previous model construction, we show the training and prediction performance using the MS-D Net combined with the relative entropy loss function and rescaled output data using empirical partition functions applied on the tKdV equations. In the numerical tests, we first consider the optimal prediction skill using the deep neural network within one updating cycle. Especially, we are more interested in capturing the non-Gaussian statistics from the network rather than the exact recovery of the single time-series which should give large difference with small perturbations due to its turbulent nature. The basic strategy is to train the model using data from the near-Gaussian solutions with a small inverse temperature θ_0 in the Gibbs measure [2]; then the optimized neural network is used to predict the more skewed model statistics for regimes with larger absolute values of θ .

For simplicity, we set the input and output data of the network in the same shape. Then the model output from one single iteration of the network is compared with the target data from the true model solution. In the structure of the neural network, the standard setup is adopted for the tKdV solutions. In each layer of the network, a kernel with minimum size $w = 3$ is taken for the convolution update. The dilation size for mixed scales changes from $s = 0$ to $s = 5$ repeatedly as the network grows in depth. Different network depths of layers L are tested for the model performance, while it is found that a moderate choice $L = 80$ is enough to produce desirable training and prediction results.

In calibrating the errors from model predictions, we propose the normalized square error between the true target \mathbf{y}^t and the network output \mathbf{y}^o

$$E(\mathbf{y}^t, \mathbf{y}^o) = \frac{\sum_{i=1}^{JN} |f(y_i^t) - f(y_i^o)|^2}{\sum_{i=1}^{JN} |f(y_i^t)|^2}, \quad [9]$$

where the subscript i represents the i -th component in the training/prediction set $\mathbf{y} \in \mathbb{R}^{J \times N}$. The function $f(y)$ acting on each component of \mathbf{y} can be used to extract the useful features to be calibrated. In the following tests, we use $f(y) = y$ to compare the original output of the data, and use the exponential scaling $f(y) = \exp(y)$ to check the prediction in positive extreme values.

Training the neural network using near-Gaussian data. In the training process for turbulent system statistics, we include the temporal and spatial correlations together in the input data by considering a short time-series of the solution. The training data is drawn from the model solutions of [1] only among the near-Gaussian regime statistics. In summary, we use the training data set in the following structure:

- The *input data* is from the ensemble solutions $u_j^{(m)}(t_n)$ of the tKdV equation. It is organized in a tensor of the shape (M, J, N) , where $M = 10000$ is the total ensemble size, $J = 32$ is the spatial discretization size, and $N = 50$ is the sampled time instants with the time step $\Delta t = 0.01$. Thus the initial data for training is given as the tKdV solutions in the time window $[0, 0.5]$.
- The *target data* for the training result is the solutions $u_j^{(m)}(T + t_n)$ of the tKdV solution. The data is organized into the tensor with the same size (M, J, N) as the input

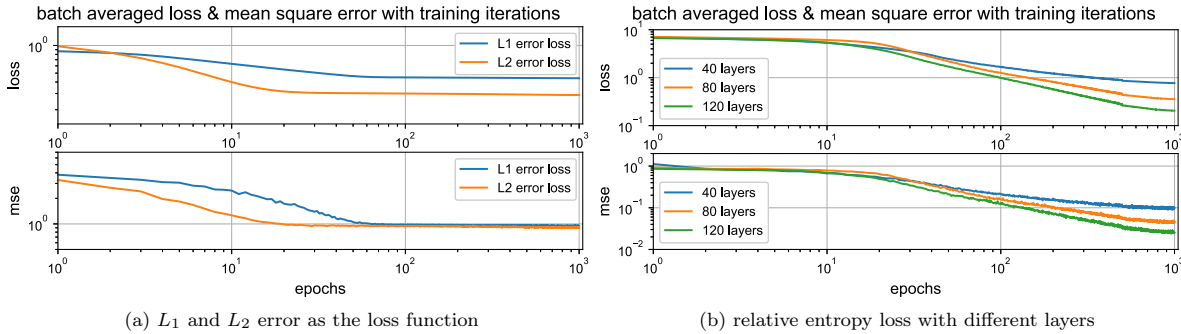


Fig. 2. Training loss function and the mean relative square error in the data using L_1 and L_2 error loss functions (left) and using the relative entropy loss function with rescaled output data (right) during the training iterations. Both networks with L_1 and L_2 loss are set to have $L = 80$ densely connected layers; and the networks with the relative entropy loss are compared using $L = 40, 80, 120$ layers.

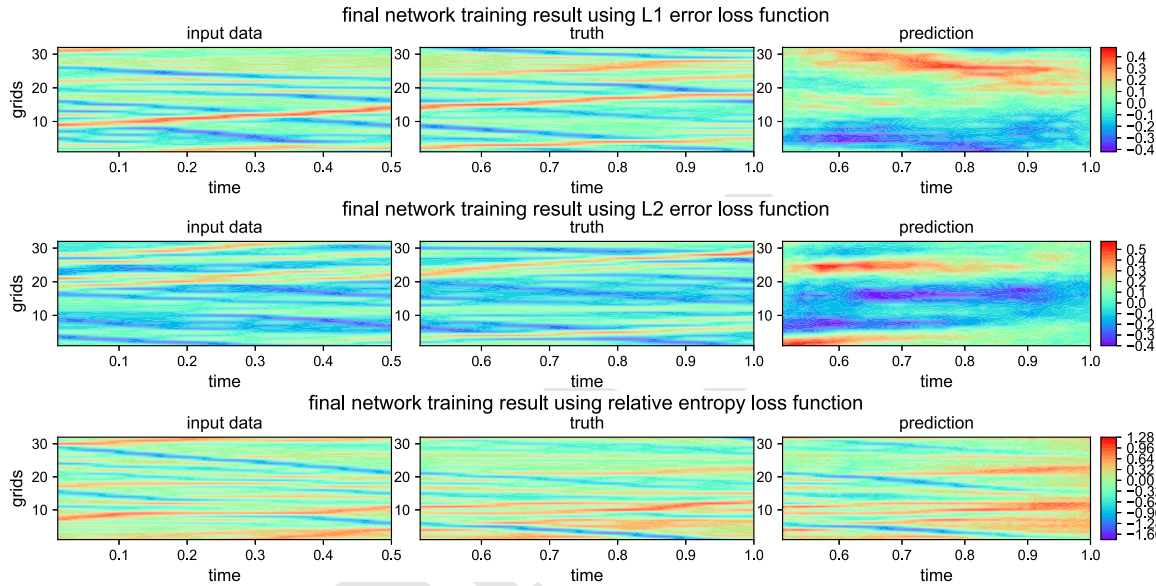


Fig. 3. One snapshot of the final training results with three different loss functions. The left panel shows the input data, the middle panel is the true target to fit, and the right panel shows the output results from the trained networks. All networks contain $L = 80$ layers in the tests.

439 data. We can consider different starting time for the target
 440 data by changing the time length T . The prediction time
 441 scale here is taken as $T = 0.5$, that is, to consider the
 442 prediction in time window $[0.5, 1]$.

- 443 • The input data (with total ensemble size 10000) is divided
 444 into 100 mini-batches with size 100 for each training group
 445 in one epoch. In total we use 1000 epochs in the entire
 446 training process. The mini-batches are randomly selected
 447 with the batch indices resampled from random numbers
 448 in each step.

449 Notice that the above prediction time length $T = 0.5$ used
 450 in the experiments is much longer beyond the decorrelation
 451 time of the tKdV states. From the direct numerical results in
 452 Figure 1, the autocorrelation function decays to near zero at
 453 $t = 0.5$, and the longest decorrelation time among the spectral
 454 modes is below 0.1.

455 First, we compare the training performance using the stan-
 456 dard L_1 and L_2 loss functions in [3] and [4] with the same
 457 MS-D Net structure. The left panel of Figure 2 shows the evo-
 458 lution of training loss functions and the mean relative square
 459 errors [9] among the training samples during the stochastic

460 gradient descent iterations. According to the loss functions
 461 under the L_1 and L_2 distances, the training appears to be
 462 effective and the error quickly dropped to smaller values in the
 463 first few steps. However, if we compare more carefully about
 464 the relative errors in the results, both cases get saturated
 465 quickly at a high error level near 1. The errors then become
 466 difficult to improve by training with larger number of itera-
 467 tions and applying deeper layers. This is because under both
 468 metrics, the model tries hard to fit the small-scale turbulent
 469 fluctuations in small values while missing the most important
 470 large-scale events in the solutions. It can be seen more clearly
 471 in Figure 3 for typical training output snapshots compared
 472 with the truth. No desirable prediction can be reached.

473 In contrast, significant improvement is achieved by switch-
 474 ing to the relative entropy loss function and adopting empirical
 475 partition functions to normalize the output data. In this case
 476 as illustrated on the right panel of Figure 2, both the relative
 477 entropy loss function and the relative error drop to very small
 478 values in the final steps of the training iterations, implying
 479 high skill of the network to produce accurate predictions in the
 480 prediction time range (Though it appears that the accuracy

Table 1. Mean and variance of the relative square errors among a test with 500 samples for the state u and the scaled state $\exp(u)$.

error		near Gaussian	mildly skewed	highly skewed
u	mean	0.2682	0.2556	0.2690
	variance	0.0039	0.0048	0.0087
$\exp(u)$	mean	0.0733	0.0764	0.0985
	variance	0.0005	0.0011	0.0060

statistics is presented. As shown in the results, the trained network displays uniform skill among all the tested samples in capturing the exact dynamical solutions in the extreme event domain unknown from the training data. By looking at the errors in the scaled data using the exponential function, the error amplitude even becomes smaller, confirming the accurate characterization of large extreme values through the network. In the typical snapshot of one sample, both the extreme values in the transporting waves and non-extreme detailed turbulent fluctuating structures are captured by the model.

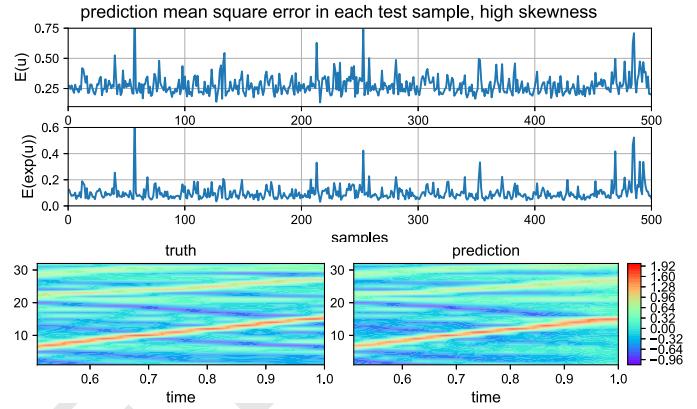


Fig. 4. Prediction in the regime with highly skewed statistics using the trained neural network with $L = 80$ layers. The upper row plots the relative square errors for the state u and the scaled state $\exp(u)$ among the 500 tested samples. The lower row shows one typical snapshot of the prediction.

Furthermore, the prediction errors as the absolute difference between the truth and the model prediction in the three statistical regimes are displayed in Figure 5. The network predictions maintain accurate with small errors in a much longer time scale than the decorrelation time $T_{\text{decorr}} < 0.1$ of the states. This gives the final confirmation of the general high skill in the deep neural network for capturing key representative features in complex dynamical systems once the essential structures are learned from the training procedure.

could be further improved by applying more iteration steps in the training, the results are already good enough after about 200 iterations). Note that we use logarithmic coordinates so the small values are emphasized. According to the last row of Figure 3 for one typical training result snapshot, both the extreme values and the small amplitude structures are captured in the model.

As a final remark, we check the proper depth needed for accurate predictions in the neural network. The right panel of Figure 2 also compares the same network under the relative entropy loss but using different numbers of layers. A deeper network clearly can further improve the prediction skill and push the final optimized error to an even lower value, with the cost of a larger computational requirement. Still from the comparison, it shows that a moderate number of layers (such as $L = 80$) is sufficient to produce accurate results with relatively low cost. By pushing the network to deeper layers with $L = 120$, the improvement in error just becomes small and may not be necessary with the additional computational cost. The last row of Figure 3 shows already the quite accurate recovery of the solution field purely learned from data.

Predicting extreme events using deep neural network. In checking the prediction skill of the optimized network, we pick the neural network with $L = 80$ densely connected layers as the standard model to test its predictions among different statistical regimes. It has been shown in the training process with a high skill in recovering the original flow structures. Next, we should confirm that the neural network has really learned the dynamical structure of the original model, instead of merely overfitting the data.

Three statistical regimes ranging from near Gaussian, mildly skewed, and highly skewed PDFs as shown in Figure 1 are taken for testing the range of prediction skill in the neural network model. An ensemble of 500 new trajectories from the tKdV solutions in different statistical regimes is used to show the robustness of the method. The mean and variance of the relative square errors [9] among the samples for the state u and the errors under the exponential function $\exp(u)$ are list in Table 1 for different statistical regimes. Uniform high accuracy in the mean with tiny variance is achieved among the vastly different regimes with distinct statistical features.

Especially, we are interested in the case with highly non-Gaussian statistics representing the frequent occurrence of extreme events. In Figure 4, the network is used to predict the flow solutions in the regime with highly skewed statistics (results for the other two cases can be found in *SI Appendix, B.2*). The extreme values are represented by high peaks of a dominant wave moving along the field. This feature is not shown at all in the training data where only near Gaussian

549 Concluding remarks

550 A new strategy using a densely connected multi-scale deep
551 neural network with relative entropy loss function for cali-
552 brating rescaled output data is proposed for the prediction of
553 extreme events and anomalous features from data. It needs
554 to be noticed that the extreme events are often represented
555 by highly skewed PDFs and have frequent occurrence in the
556 turbulent field (3, 8, 12) in contrast to the other situation
557 of isolated rare events which can be studied with machine
558 learning models (9). The prediction skill of the optimized deep
559 neural network is tested on the truncated KdV equation, where
560 different Gibbs states create a wide range of statistics from
561 near Gaussian to highly skewed distributions. By adopting the
562 densely connected and multi-scale structures, the deep neural
563 network is easy to train with standard model setup and fewer
564 model hyperparameters.

565 Using training data only drawn from the near-Gaussian
566 regime of the dynamical model, the deep neural network dis-
567 plays high skill in learning the essential dynamical structures
568 of the complex system and provides uniformly accurate pre-
569 diction among a wide range of different regimes with distinct
570 statistics. The network also shows robustness among tests in a
571 large ensemble of samples. The robust performance in the test
572 model implies the potential of more general applications using
573 the neural network framework for the prediction of extreme
574 events and important statistical features in a wider group of
575 more realistic high-dimensional turbulent dynamical systems.

576 **ACKNOWLEDGMENTS.** This research of A.J.M. is partially
577 supported by the Office of Naval Research N00014-19-S-B001. D.Q.
578 is supported as a postdoctoral fellow on the grant.

- 579 1. Bolles CT, Speer K, Moore MNJ (2019) Anomalous wave statistics induced by abrupt depth
580 change. *Physical Review Fluids* 4(1):011801.
- 581 2. Dematteis G, Grafke T, Vanden-Eijnden E (2018) Rogue waves and large deviations in deep
582 sea. *Proceedings of the National Academy of Sciences* 115(5):855–860.
- 583 3. Qi D, Majda AJ (2018) Predicting extreme events for passive scalar turbulence in two-layer
584 baroclinic flows through reduced-order stochastic models. *Communications in Mathematical*
585 *Sciences* 16:17–51.
- 586 4. Mohamad MA, Sapsis TP (2018) Sequential sampling strategy for extreme event statis-
587 tics in nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*
588 115(44):11138–11143.
- 589 5. Qi D, Majda AJ (2016) Predicting fat-tailed intermittent probability distributions in passive
590 scalar turbulence with imperfect models through empirical information theory. *Communica-*
591 *tions in Mathematical Sciences* 14(6):1687–1722.
- 592 6. Majda AJ, Tong XT (2015) Intermittency in turbulent diffusion models with a mean gradient.
593 *Nonlinearity* 28(11):4171.
- 594 7. Viotti C, Dias F (2014) Extreme waves induced by strong depth transitions: Fully nonlinear
595 results. *Physics of Fluids* 26(5):051705.
- 596 8. Cai D, Majda AJ, McLaughlin DW, Tabak EG (2001) Dispersive wave turbulence in one di-
597 mension. *Physica D: Nonlinear Phenomena* 152:551–572.
- 598 9. Guth S, Sapsis TP (2019) Machine learning predictors of extreme events occurring in complex
599 dynamical systems. *Entropy* 21(10).
- 600 10. Cousins W, Sapsis TP (2015) Unsteady evolution of localized unidirectional deep-water wave
601 groups. *Physical Review E* 91(6):063204.
- 602 11. Majda AJ, Lee Y (2014) Conceptual dynamical models for turbulence. *Proceedings of the*
603 *National Academy of Sciences* 111(18):6548–6553.
- 604 12. Grooms I, Majda AJ (2014) Stochastic superparameterization in a one-dimensional model for
605 wave turbulence. *Communications in Mathematical Sciences* 12(3):509–525.
- 606 13. Chen N, Majda AJ (2017) Beating the curse of dimension with accurate statistics for the
607 Fokker–Planck equation in complex turbulent systems. *Proceedings of the National Academy*
608 *of Sciences* 114(49):12864–12869.
- 609 14. Majda AJ, Moore MNJ, Qi D (2019) Statistical dynamical model to predict extreme events and
610 anomalous features in shallow water waves with abrupt depth change. *Proceedings of the*
611 *National Academy of Sciences* 116(10):3982–3987.
- 612 15. Majda AJ, Qi D (2019) Statistical phase transitions and extreme events in shallow water
613 waves with an abrupt depth change. *submitted to Journal of Statistical Physics*.
- 614 16. Schmidhuber J (2015) Deep learning in neural networks: An overview. *Neural networks*
615 61:85–117.
- 616 17. LeCun Y, Bengio Y, Hinton G (2015) Deep learning. *Nature* 521(7553):436.
- 617 18. Jordan MI, Mitchell TM (2015) Machine learning: Trends, perspectives, and prospects. *Sci-*
618 *ence* 349(6245):255–260.
- 619 19. Goodfellow I, Bengio Y, Courville A (2016) *Deep learning*. (MIT press).

20. Pathak J, Hunt B, Girvan M, Lu Z, Ott E (2018) Model-free prediction of large spatiotempo-
620 rally chaotic systems from data: A reservoir computing approach. *Physical review letters*
621 120(2):024102.
- 622 21. Raissi M, Karniadakis GE (2018) Hidden physics models: Machine learning of nonlinear
623 partial differential equations. *Journal of Computational Physics* 357:125–141.
- 624 22. Bolton T, Zanna L (2019) Applications of deep learning to ocean data inference and subgrid
625 parameterization. *Journal of Advances in Modeling Earth Systems* 11(1):376–399.
- 626 23. Weyn JA, Durran DR, Caruana R (2019) Can machines learn to predict weather? using deep
627 learning to predict gridded 500-hpa geopotential height from historical weather data. *Journal*
628 *of Advances in Modeling Earth Systems*.
- 629 24. Han J, Jentzen A, E W (2018) Solving high-dimensional partial differential equations using
630 deep learning. *Proceedings of the National Academy of Sciences* 115(34):8505–8510.
- 631 25. E W, Han J, Jentzen A (2017) Deep learning-based numerical methods for high-dimensional
632 parabolic partial differential equations and backward stochastic differential equations. *Com-*
633 *munications in Mathematics and Statistics* 5(4):349–380.
- 634 26. Brenowitz ND, Bretherton CS (2018) Prognostic validation of a neural network unified physics
635 parameterization. *Geophysical Research Letters* 45(12):6289–6298.
- 636 27. Pelt DM, Sethian JA (2018) A mixed-scale dense convolutional neural network for image
637 analysis. *Proceedings of the National Academy of Sciences* 115(2):254–259.
- 638 28. Majda AJ, Wang X (2006) *Nonlinear dynamics and statistical theories for basic geophysical*
639 *flows*. (Cambridge University Press).
- 640