

1 **Bringing statistics to storylines: rare event sampling**
2 **for sudden, transient extreme events**

3 **Justin Finkel¹, Paul A. O’Gorman¹**

4 ¹Department of Earth, Atmospheric, and Planetary Sciences, Massachusetts Institute of Technology

5 **Key Points:**

- 6 • Rare event algorithms may help address the challenge of simulating extreme weather
7 events and quantifying their probability.
8 • When the event of interest is sudden and transient, perturbed ensembles diver-
9 sify too slowly for standard rare event algorithms to work.
10 • Using the Lorenz-96 model as a prototype for midlatitude weather, we use early
11 perturbation and a rejection step to gain a speedup.

Corresponding author: Justin Finkel, ju26596@mit.edu

Abstract

A leading goal for climate science and weather risk management is to accurately model both the physics and statistics of extreme events. These two goals are fundamentally at odds: the higher a computational model’s resolution, the more expensive are the ensembles needed to capture accurate statistics in the tail of the distribution. Here, we focus on events that are localized in space and time, such as heavy precipitation events, which can start suddenly and decay rapidly. We advance a method for sampling such events more efficiently than straightforward climate model simulation. Our method combines elements of two recent approaches: adaptive multilevel splitting (AMS), a rare event algorithm that generates rigorous statistics at reduced cost, but that does not work well for sudden, transient extreme events; and “ensemble boosting” which generates physically plausible storylines of these events but not their statistics. We modify AMS by splitting trajectories well in advance of the event’s onset following the approach of ensemble boosting, and this is shown to be critical for amplifying and diversifying simulated events in tests with the Lorenz-96 model. Early splitting requires a rejection step that reduces efficiency, but nevertheless we demonstrate improved sampling of extreme local events by a factor of order 10 relative to direct sampling in Lorenz-96. Our work makes progress on the challenge posed by fast dynamical timescales for rare event sampling, and it draws connections with existing methods in reliability engineering which, we believe, can be further exploited for weather risk assessment.

Plain Language Summary

What is the strongest rainstorm that we can expect in a given thousand-year period? To augment the available ~ 100 years of historical data and to account for climate change, computer simulations are a useful, but expensive, tool to answer such questions. A model must run for many millennia to deliver an answer with statistical confidence. *Rare event algorithms* provide a promising alternative simulation protocol, in which an ensemble of short simulations is biased to produce more extreme events and reweighting is used to correct for the bias when calculating statistics. However, a classical rare event algorithm fails when the events of interest are short and “bursty” (like heavy rainstorms) instead of long and slow-moving (like anomalously hot summers). We modify the rare event algorithm to make it amenable to precipitation-like events in an idealized dynamical system with chaotic traveling waves.

1 Introduction

In climate modeling, high spatial resolution is important for realistically representing localized extreme weather events like cyclones producing extreme precipitation and winds (O’Brien et al., 2016; van der Wiel et al., 2016). But given finite computational resources, high resolution has to be traded off with the need for ensembles of models and simulations to deal with uncertainty related to model physics, parameters, initial conditions and boundary conditions including emissions scenarios. Extreme events are particularly challenging because they occur infrequently, and hence need large ensemble sizes to have their small probabilities accurately quantified. The conflict for computational resources therefore comes to a head in the study of extreme events.

A variety of shortcuts have developed in the past century to alleviate this conflict. Leading statistical approaches include extreme value theory (EVT; Coles, 2001) and large deviation theory (Touchette, 2009), which respectively describe the behavior of *maxima* and anomalously large *running means* in random processes. In principle, we can use these theories to fit a parametric family to limited data and then extrapolate to even longer return periods. EVT has become an important tool in risk assessment and climate change attribution (Kharin et al., 2007; Naveau et al., 2020), while large deviation theory succinctly encodes the severity of long-lasting, large-area events such as persistent heat waves

(Gálfi et al., 2021). Statistical theories help make the most of a fixed dataset, but parameter estimation can be unstable given the restrictive underlying assumptions and the limited datasets available (W. K. Huang et al., 2016; Gálfi et al., 2017). For example, EVT only holds in the limit of large blocks of data or high thresholds for extremity, which directly conflicts with the requirement of many samples for low-variance parameter estimation. Moreover, statistical theories don't provide spatio-temporal resolved extreme events (e.g., the spatial field of rainfall and other fields on the day of an extreme event) which are needed to drive impact models.

Statistical or dynamical downscaling is another way to address the problem of extremes by reducing the computational cost of obtaining high-resolution output from long simulations or large ensembles (X. Huang et al., 2020; Lee et al., 2020; Emanuel, 2021; Saha & Ravela, 2022; Krouma et al., 2022). Downscaling nevertheless has some drawbacks. Dynamical downscaling using regional climate models faces the challenge of correctly forcing a regional model with output from a different global model, and the regional model inherits errors in large-scale fields from the global model (Adachi & Tomita, 2020), while statistical downscaling assumptions can create systematic errors (Schmidli et al., 2007) and may not generalize to different climates.

The focus of this paper is *rare event sampling*, which is a strategy for allocating more of the computational effort towards rare events, and less effort towards the long intervening periods of comparatively mild behavior. This is usually achieved by *splitting* methods, which consist of three steps repeated in a cycle: (1) run an ensemble of simulations forward, (2) identify the ensemble members making the most progress towards the extreme event, and (3) clone these most-promising ensemble members (applying small perturbations) while discarding the less-promising members, resulting in a new ensemble that is more prone to extremes than was the original ensemble. With repeated rounds of splitting, one can populate the tail of the probability distribution more fully, while neglecting the more typical behavior of lesser interest. Crucially, in statistical analysis of the ensemble, one must compensate for the bias by weighting each clone with a factor less than one, relying on the *importance sampling* formalism. See Bucklew (2004) for an introduction to rare event sampling.

This generic procedure has many possible variants, which have been developed largely in the fields of physics (Kahn & Harris, 1951; Giardinà et al., 2006), chemistry (Kästner, 2011; Zuckerman & Chong, 2017), and reliability engineering (Au & Beck, 2001), but have recently started to make an impact on Earth and planetary sciences. For example, extreme European heat waves were sampled by Ragone et al. (2018) and Ragone and Bouchet (2021) with genealogical particle analysis (GPA), and by Yiou and Jezequel (2020) with empirical importance sampling. Wouters et al. (2023) sampled extreme European seasonal precipitation accumulations, also using GPA. Webber et al. (2019) developed a quantile-based variant of GPA to sample more extreme versions of tropical cyclones. Planetary science applications include jet nucleation (Bouchet et al., 2019) and orbit destabilization (Abbot et al., 2021). For studies of climate, rare event sampling can be applied to global models or paired with the dynamical and statistical downscaling approaches mentioned earlier.

We have elected to use a particular rare event algorithm called *adaptive multilevel splitting* (AMS), which was first established by Cérou and Guyader (2007) and is similar to the earlier RESTART algorithm (Villén-Altamirano et al., 1991). Lestang et al. (2018) successfully applied AMS to the Ornstein-Uhlenbeck process, while Lucente, Roland, et al. (2022) and Baars et al. (2021) used AMS to study regime transitions in idealized climate models. AMS has also been usefully employed in other diverse fields such as molecular dynamics and air traffic control (see Cérou et al. (2019) for a recent review). The distinguishing feature of AMS is that it operates on the level of full trajectories over a fixed time horizon, and applies the small perturbation to trajectories at the instant that they first cross a threshold of extremity. The “child” trajectory is identical to its par-

115 ent up until this time, whereas it diverges from its parent afterward to give a new re-
 116 realization of the extreme event. All ensemble members failing to cross the threshold are
 117 discarded, and the threshold is then raised for repeated rounds of splitting and killing.

118 A related approach, “ensemble boosting”, is a novel technique for generating “sto-
 119 rylines” of unprecedented climate extremes (Gessner et al., 2021; Gessner, 2022). In this
 120 approach, one identifies several extreme events from a long climate simulation, perturbs
 121 the antecedent conditions (1-3 weeks ahead of time), and re-simulates the event to gen-
 122 erate alternative realities, which sometimes turn out even more extreme. While similar
 123 to splitting methods, ensemble boosting does not explicitly quantify statistics. As ex-
 124 plained below, a major goal of this paper is to combine the benefits of ensemble boost-
 125 ing with that of rare event algorithms, in particular AMS.

126 Given the successes in using rare event sampling discussed above, it is desirable to
 127 also use it to sample shorter-term extreme weather events, such as daily precipitation
 128 extremes, which have large societal impacts in the current climate (Wright et al., 2021;
 129 Thompson et al., 2017) and are expected to intensify under climate change (O’Gorman,
 130 2015; Pfahl et al., 2017; Tandon et al., 2018; Myhre et al., 2019). However, heavy pre-
 131 cipitation events (or high wind events) have some dynamical characteristics that distin-
 132 guish them from the previous applications and pose challenges to existing rare event al-
 133 gorithms. Unlike continental-scale, seasonally averaged anomalies studied previously (Ragone
 134 et al., 2018; Wouters et al., 2023), heavy precipitation events of interest are often sud-
 135 den, transient, and relatively small-scale. Their timescale at a particular location is of-
 136 ten limited by the propagation of the dynamical feature causing the precipitation such
 137 as cyclones and fronts (Dwyer & O’Gorman, 2017). The strategy used in Ragone et al.
 138 (2018) and Wouters et al. (2023) relies on some slow-moving notion of *progress* towards
 139 the extreme event, naturally given by the integrated temperature anomaly itself when
 140 targeting extreme seasonal average temperatures, in order to decide which simulations
 141 to clone or kill. In the precipitation study of Wouters et al. (2023), the extreme event
 142 is again a seasonal total, for which a mid-seasonal total is a reasonable measure of progress.
 143 But for individual precipitation events, if one uses precipitation itself to measure progress
 144 towards the event, and applies perturbations to a simulation when precipitation picks
 145 up, it is too late for these perturbations to take effect by the time of maximum precip-
 146 itation. The event simply comes and goes faster than perturbed simulations diverge. Lestang
 147 et al. (2018) found a similar pathology with AMS when sampling extreme pressure fluctu-
 148 ations on a body embedded in a turbulent channel flow. There, the extreme events were
 149 caused by vortices sweeping past the body, roughly analogous to cyclones sweeping past
 150 a location on Earth, and the rapidity of the fluctuation crippled the effectiveness of the
 151 standard splitting strategy.

152 To isolate and solve the problem of applying rare event algorithms to sudden, tran-
 153 sient extremes, we postpone the specific application to precipitation and first descend
 154 the model hierarchy to the Lorenz-96 model (Lorenz, 1996), a spatiotemporal chaotic
 155 system often used as a toy model for the atmosphere. The model produces extreme events
 156 posing the same algorithmic challenges as precipitation extremes: intermittent, short-
 157 lived bursts carried by traveling waves with unpredictable amplitudes. It has been used
 158 in numerous past studies of extreme event statistics and predictability (Sterk & van Kekem,
 159 2017; Qi & Majda, 2016; Hu et al., 2019). With this cheap but behaviorally rich model,
 160 we have developed a simple modification to AMS, drawing inspiration from ensemble boost-
 161 ing by simply applying a split in advance of the event’s onset by some advance split time
 162 δ —hence, “trying early” AMS (TEAMS). To make this statistically rigorous, a rejection
 163 step is necessary, which comes at an efficiency cost, but still enables moderate speedups
 164 of ~ 10 relative to direct sampling. Fig. 1 displays a schematic diagram for TEAMS, which
 165 will be elaborated in section 3. In fact, TEAMS is a repurposing of a more general method
 166 called *subset simulation* (Au & Beck, 2001) from structural reliability engineering, a field
 167 whose sophisticated rare event algorithms could benefit the climate risk community.

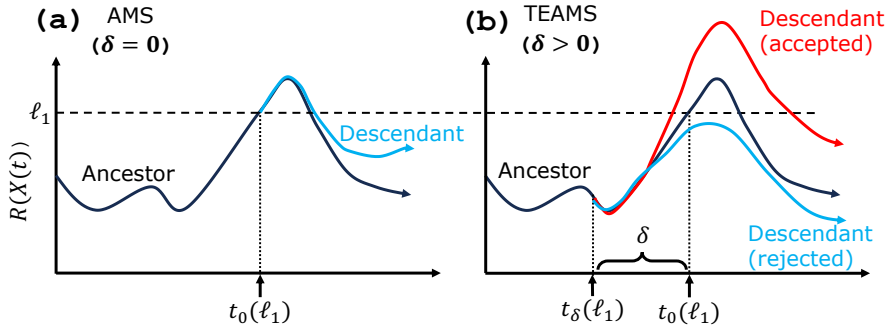


Figure 1. Schematic of the splitting step in (a) AMS and (b) TEAMS. Black curves represent an initial ensemble member, or ancestor, which exceeds the first level ℓ_1 and has been selected for cloning in the first round. In AMS, the perturbation is applied at the instant $t_0(\ell_1)$ when the ancestor first exceeds ℓ_1 , resulting in a descendant trajectory (blue) which essentially replicates the extreme event because the separation timescale is longer than the event itself. On the other hand, in TEAMS (right) we apply the perturbation in advance, by some margin $\delta > 0$. This can sometimes result in rejection (blue descendant), i.e., failure to cross ℓ_1 . However, when a descendant is accepted (red) it will be more distinct from the ancestor than the corresponding descendant in AMS and have the potential to reach a substantially higher peak value.

168 This paper is organized as follows. In section 2, we present a stochastically forced
 169 Lorenz-96 model and the behavior of its extreme events as a function of stochastic forc-
 170 ing strength. In section 3, we first introduce the general framework of subset simulation.
 171 In section 3.1, we specialize to AMS, and in section 3.2 we show that AMS fails in the
 172 low-noise forcing regime, which is often most relevant for weather and climate models.
 173 In section 3.3, we modify AMS to use a “trying early” step with rejection sampling and
 174 recover a substantial speedup. In section 4, we further explore the relationship between
 175 the advance splitting time—a key algorithmic parameter—and classical notions of pre-
 176 dictability timescales. Finally, in section 5 we point out directions for further develop-
 177 ment.

178 2 Lorenz-96: a customizable spatiotemporal chaotic system

179 Lorenz (1996) introduced a simple dynamical system (L96 hereafter) meant to cap-
 180 ture some crucial aspects of atmospheric dynamics. The model state space consists of
 181 $K (\geq 4)$ variables $\{x_k\}_{k=1}^K$ arranged on a one-dimensional periodic lattice, each k rep-
 182 resenting a longitude sector on Earth. x_k represents a generic atmospheric variable like
 183 wind speed or vorticity and evolves according to the coupled equations

$$\frac{dx_k}{dt} = ax_{k-1}(x_{k+1} - x_{k-2}) - x_k + \mathcal{F}_k, \quad k = 0, \dots, K - 1, \quad (1)$$

184 where x_{k+K} is identified with x_k . The quadratic terms on the right-hand side represent
 185 advection, like the quadratic nonlinearity in the material derivative of the Navier-Stokes
 186 equations, which on its own conserves “energy” $\frac{1}{2} \sum_k x_k^2$. The linear term $-x_k$ repre-
 187 sents damping due to friction, and the additive term \mathcal{F}_k represents external forcing, like
 188 a meridional insolation gradient. The latter two terms destroy exact energy conserva-
 189 tion, but balance out in a time-averaged sense to make for a statistically steady state.
 190 Lorenz (1996) introduced the above model with \mathcal{F}_k constant in k and also a version in
 191 which \mathcal{F}_k is a “subgrid-scale forcing” that is a function of an additional tier of dynam-

Table 1. Physical parameters for Lorenz-96 system (upper section), and algorithmic parameters for the TEAMS algorithm (lower section).

Symbol	Explanation	Value or range
K	Number of longitude sites	40
a	Strength of advection term	$\{1, 0\}$ (mostly 1)
F_0	Constant background forcing	6
m	Wavenumber for stochastic forcing	$\{1, 4, 7, 10\}$ (mostly 4)
F_m	Strength of stochastic forcing at wavenumber m	$\{3, 1, 0.5, 0.25, 0\}$
N	Number of initial ensemble members	128
κ	Number of members to kill each round	1
J	Number of rounds of splitting	896
T	Time horizon	6
δ	Advance split time	$[0, 2]$

192 ical variables representing finer scales, and this version has proven useful for testing stochastic
 193 parameterization schemes (e.g., Wilks, 2005; Hu et al., 2019; Gagne II et al., 2020).
 194 Here, we also allow \mathcal{F}_k to vary stochastically with longitude (k) and time:

$$\mathcal{F}_k = F_0 + F_m \left[\eta_1 \cos \left(\frac{2\pi mk}{K} \right) + \eta_2 \sin \left(\frac{2\pi mk}{K} \right) \right] \quad (2)$$

195 where $\eta_{1,2}$ are independent Gaussian white-noise processes, and m is an integer wavenum-
 196 ber. Formally, Eq. (2) renders Eq. (1) a diffusion process, using the Itô convention for
 197 stochastic integrals (Pavliotis, 2014). This simple stochastic forcing is analagous to a stochas-
 198 tic parameterization in a weather or climate model, and in the AMS framework it allows
 199 us to easily generate new ensemble members by splitting an existing ensemble member
 200 at a certain time. We verify below that for weak amplitudes the stochastic forcing does
 201 not appreciably alter model statistics.

202 The parameters used here are summarized in the upper section of Table 1. We set
 203 $K = 40$, following Lorenz and Emanuel (1998). We fix the constant part of the forc-
 204 ing to be $F_0 = 6.0$, which is sufficient for weak turbulence (a larger value would be needed
 205 with smaller K). We choose the stochastic forcing wavenumber as $m = 4$ because that
 206 empirically seems to drive ensemble members apart slightly faster than very small or large
 207 wavenumbers (see section 4.2). Indeed the stochastically perturbed parameterization ten-
 208 dencies (SPPT) method developed at ECMWF uses noise that is spatially correlated at
 209 a $\sim 10^\circ$ length scale (Buizza et al., 1999; Palmer et al., 2009). The amplitude of $F_m (=$
 210 $F_4)$ will be explored systematically below. One further parameter, the coefficient a , de-
 211 termines the strength of the advection term. $a = 1$ is standard for L96, while $a = 0$
 212 gives an array of correlated Ornstein-Uhlenbeck (OU) processes (Pavliotis, 2014). Re-
 213 taining the OU process as a special case of L96 is useful to provide a reference case on
 214 which existing rare event splitting algorithms excel. Results for $a = 0$ are shown in sup-
 215plementary Figs. S1 and S2, and all other results presented are for $a = 1$.

216 Fig. 2 displays short numerical integrations of L96 with four different parameter
 217 choices. We used the Euler-Maruyama method with a timestep of 0.001 to integrate Eq. (1),
 218 saving out every 0.05 time units. For comparison, Lorenz and Emanuel (1998) interpret
 219 a single time unit as 5 days. The left column shows single-site variables $x_0(t)$ for each
 220 parameter set, while the right column shows corresponding Hovmöller diagrams. In the
 221 standard deterministic system $F_4 = 0$ in the top row, $x_0(t)$ fluctuates with a semi-regular
 222 period of ~ 2 time units (10 “days”) but with irregular amplitudes, the largest of which
 223 are precisely the extreme events we choose to study here. The Hovmöller diagram re-

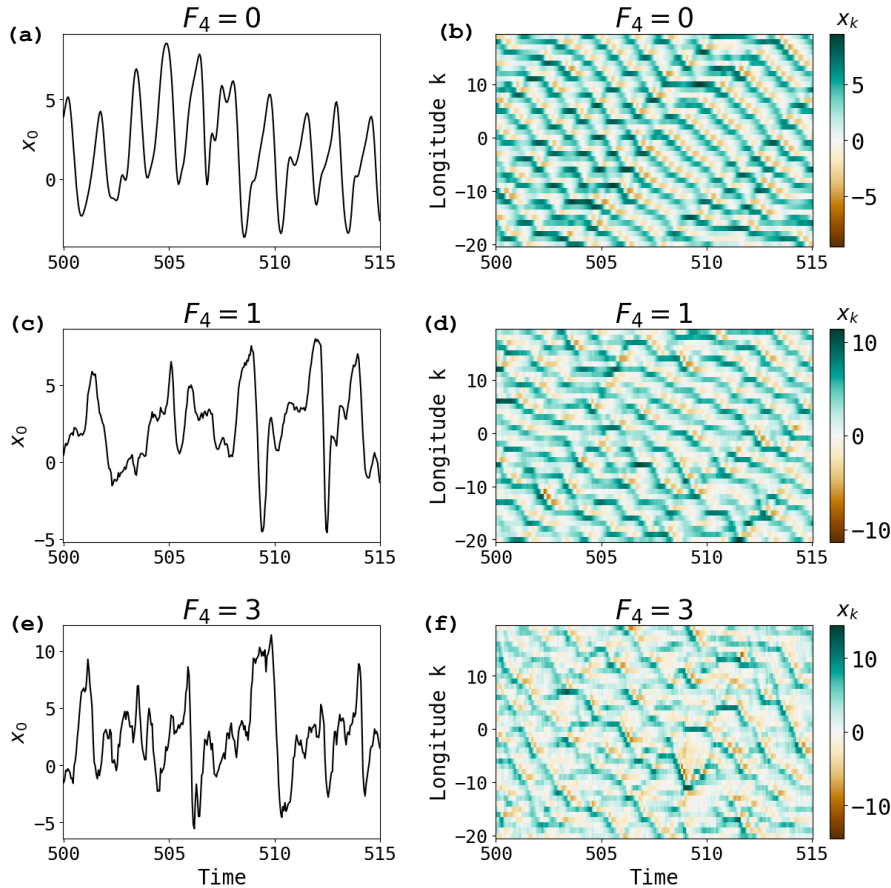


Figure 2. Time evolution of the L96 model expressed as timeseries of $x_0(t)$ (left column) and Hovmöller diagrams (right column) with three different levels of stochastic forcing. (a,b) have $F_4 = 0$ (the deterministic system); (c,d) have $F_4 = 1$ (moderate forcing); (e,f) have $F_4 = 3$ (strong forcing).

224 veals these fluctuations to arise from a field of traveling waves, with roughly eight peaks
 225 and troughs moving with negative (“westward”) phase velocity. The waves experience
 226 intermittent disturbances, sometimes getting stuck in place for several turnover times
 227 and setting up favorable conditions for extreme events. Globally, these stagnations man-
 228 ifest as kinks that propagate in the positive (“eastward”) direction. This is reminiscent
 229 of atmospheric Rossby waves, whose phase and group velocities have opposite signs (up
 230 to a Doppler shift due to the mean flow) (Lorenz & Emanuel, 1998). Thus, we can loosely
 231 think of the peaks and troughs as being like highs and lows in the midlatitude atmosphere.

232 Fig. 2 rows 2 and 3 show analogous pictures for moderate ($F_4 = 1$) and strong
 233 ($F_4 = 3$) stochastic forcing, respectively. As noise increases the traveling waves tran-
 234 sition from unidirectional to zigzagging. The timeseries become more jagged and more
 235 liable to take large excursions from their mean and hover there for longer durations.

236 Fig. 3a overlays PDFs of the single-site value (x_0) for all these parameter regimes,
 237 plus two more: $F_4 = 0.5$ and 0.25 . Reducing the noise roughly preserves the mode but
 238 shrinks the tails. The PDF appears basically converged for $F_4 \leq 0.5$. Fig. 3b confirms
 239 this is true even in the far tail, with a log-transformed plot of return level vs. return time
 240 for x_0^2 . The limiting case $F_4 = 0$ has a bounded tail, which is easy to see with an en-

241 ergy argument (see also Qi and Majda (2016)): defining $\bar{x} = \frac{1}{K} \sum_{k=1}^K x_k$, the energy
242 $E = \frac{1}{2} \sum_k x_k^2$ evolves as $\frac{dE}{dt} = -2E + KF\bar{x}$. Since $|\bar{x}| \leq \sqrt{x^2} = \sqrt{2E/K}$ by the
243 Cauchy-Schwarz inequality, the first term dominates for E larger than some critical E_0 ,
244 which must therefore bound the steady-state distribution's tail. However, E_0 would in-
245 crease with K , i.e., higher-dimensional systems can in principle support heavier tails (e.g.
246 Lucarini et al., 2016, ch. 4 discusses general relationships between the shape paramete-
247 r and the attractor dimension). This is part of our motivation to set K relatively large.

248 The return level vs. return period plot (as in Fig. 3b) will be used throughout the
249 paper, and we calculate it using the ‘‘modified block maximum’’ method of Lestang et
250 al. (2018). For a fixed *return level* ℓ , the *return period* $\tau(\ell)$ is defined as the mean (over
251 initial conditions and noise realizations) of the waiting time until an exceedance occurs:
252 $\tau(\ell) = \mathbb{E}[\min\{t : R(x(t)) > \ell\}]$, where R is some observable of interest for the dy-
253 namical system. We take $R(x) = x_0^2$, the local energy (times two) at longitude $k = 0$.
254 Lestang et al. (2018) approximates the exceedance times by a Poisson process for high
255 ℓ to give

$$\tau(\ell) = -\frac{T}{\log[1 - p_T(\ell)]}. \quad (3)$$

256 where $p_T(\ell)$ is the probability of at least one exceedance in a fixed time T . $p_T(\ell)$ can
257 be estimated from any collection of length- T blocks of data—*either from a single con-*
258 *tinuous timeseries or not*. This is very useful because rare event splitting algorithms gen-
259 erate branching trees of short trajectories, from which we can estimate block-wise ex-
260 ceedances but not return times directly.

261 To produce Fig. 3b, we started with simulations of length 1.28×10^6 (after dis-
262 carding the first 50 for spinup), split them into B blocks of length $T = 6$, and measure
263 the maxima M_1, \dots, M_B of x_0^2 over each block. Letting $M_{(b)}$ denote the b th largest block
264 maximum, we use the empirical (complementary) CDF estimator, $\hat{p}_T(M_{(b)}) = b/B$. Hence,
265 the return curve should interpolate the ordered pairs $(\tau_b, \ell_b) = (-\frac{T}{\log(1-b/B)}, M_{(b)})$. Be-
266 cause it is common to think of ℓ as a function of τ , and to consider logarithmically spaced
267 return periods, we linearly interpolate $M_{(b)}$ over $\log \tau_B$ to get a curve $\hat{\ell}(\tau)$. We bootstrap
268 to estimate uncertainty, resampling the blocks $1, \dots, B$ with replacement and repeating
269 the above procedure 5000 times. Shading indicates the basic bootstrap 95% confidence
270 interval (Wasserman, 2004), meaning $\hat{\ell}(\tau) + (\hat{\ell}(\tau) - \ell_{0.975}^*(\tau), \hat{\ell}(\tau) - \ell_{0.025}^*(\tau))$, where ℓ_α^*
271 denotes the α th quantile of the bootstrap distribution of ℓ for each τ . Note that when
272 $\ell_{0.025}^*(\tau)$ is much less than $\hat{\ell}(\tau)$, we get a very large *upper* bound on the confidence in-
273 terval, because it suggests via the basic bootstrap philosophy that $\hat{\ell}(\tau)$ could be very much
274 less than the true parameter $\ell(\tau)$. The lowest-noise curves are close to within uncertainty
275 even in the far tails, demonstrating the convergence of extreme value statistics for $F_4 \leq$
276 0.5. This confirms that stochastic forcing, when sufficiently weak, does not alter the sys-
277 tem's statistics very much, which allows us to approximate the deterministic system's
278 rare events while remaining within the AMS framework which relies on explicit random-
279 ness.

280 The longest return period estimable by this method of ‘‘direct numerical simula-
281 tion’’ (DNS) is $\sim 8 \times 10^5$, the simulation's length. Rare event algorithms can sample
282 physical realizations of extreme events at long return periods $\tau(\ell)$ with much less com-
283 putation time than $\tau(\ell)$, but have not yet been applied to local events in L96 with weak
284 stochastic forcing. Wouters and Bouchet (2016) did apply rare event algorithms to L96,
285 but their system parameters differed substantially from ours, with $F_0 = 256$ giving a
286 much more turbulent regime reminiscent of a stochastic process. Moreover, their target
287 quantity of interest was a globally averaged energy, whereas we target local energy at
288 one longitude as a closer analogue to extreme precipitation or winds hitting a particu-
289 lar location.

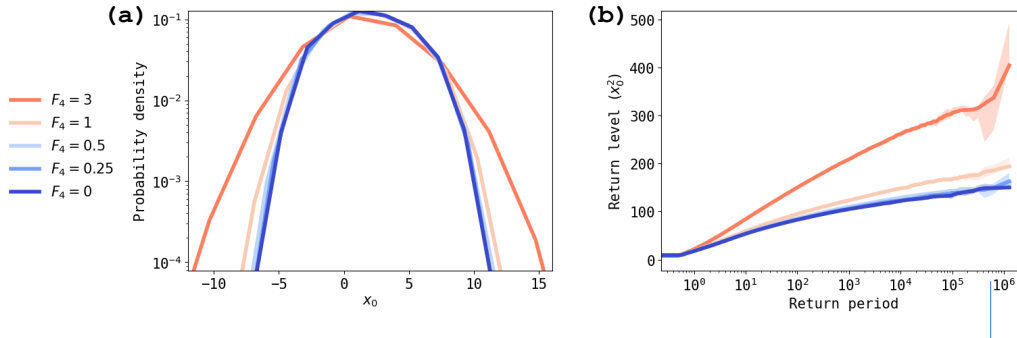


Figure 3. Steady-state statistics of the L96 model as a function of noise strength, calculated from a long simulation of length 1.28×10^6 . (a) Histograms of the model variable at one site (x_0) and (b) return level vs. return period for (twice) the local energy x_0^2 . Shading in (b) represents 95% bootstrapped confidence intervals from the modified block maximum method. See text for details.

290 The parameters a and F_4 allow us to test the performance of AMS for a range of
 291 problems, from systems on which AMS performs well to more difficult systems akin to
 292 the extreme local precipitation problem. $a = 0$ (the OU process) is an easy setting for
 293 AMS; $a = 1$ with large noise F_4 is harder, but still doable because of the dominance
 294 of noise. Shrinking F_4 further, towards the system of actual interest, gradually renders
 295 standard AMS ineffective and leads us to a modified version of the algorithm called TEAMS
 296 that allows for early splitting. The next sections present the basic algorithm and its mod-
 297 ification along this parameter path.

298 3 Subset simulation

299 TEAMS (and the special case AMS) may be viewed as a version of *subset simu-*
 300 *lation* (SS), which we use to frame our overall approach, and which we believe has con-
 301 siderable potential for application to climate problems. SS was introduced in Au and Beck
 302 (2001) and has been most widely used in structural reliability engineering (X. Huang et
 303 al., 2016). For a short pedagogical introduction, see Zuev (2015). The description be-
 304 low will introduce several tunable algorithmic parameters, which are summarized in the
 305 lower section of Table 1.

306 The goal is to estimate the probability that a random variable x from a distribu-
 307 tion ρ gives rise to large values of some quantity of interest $S(x)$,

$$p(\ell) = \int \mathbb{I}\{S(x) > \ell\} \rho(x) dx = \mathbb{E}_\rho[\mathbb{I}\{S(X) > \ell\}], \quad (4)$$

308 given only the ability to draw samples $X_1, X_2, \dots \sim \rho$. $\mathbb{I}\{\cdot\}$ denotes the indicator func-
 309 tion: one if the argument is true, zero if false. For us, each $X_i = \{X_i(t) : 0 \leq t \leq T\}$
 310 is a length- T trajectories of L96 (with stochastic forcing); the score function is a max-
 311 imum over the interval, $S(X) = \max_{0 \leq t < T} R(X(t))$; and $\rho(x)$ is the distribution over
 312 trajectories of length T induced by the stochastically forced L96 system. In structural
 313 engineering, X might be the state of a building or dam, with $\rho(x)$ induced by a proba-
 314 bility distribution over external stresses like wind, earthquakes, or rainfall, while $S(x)$
 315 would measure the proximity to failure. Because the probabilities of interest are very small,
 316 a set of independent samples $\{X_n\}_{n=1}^N$ from ρ will usually have few if any exceedances,
 317 making the “vanilla” Monte Carlo estimate of $p(\ell)$ (the fraction of exceedances) subject
 318 to high relative uncertainty. The ratio of the estimator’s variance to its mean is approx-
 319 imately $1/\sqrt{Np(\ell)}$ (Zuev, 2015). If we want to aim for a tenfold-longer return period

320 with the same uncertainty, we need to generate tenfold more samples. Worse, to reduce
 321 uncertainty tenfold we would need one hundredfold more samples, which may be unten-
 322 able.

323 SS breaks down this task into a sequence of easier tasks by setting up a series of
 324 intermediate levels $\ell_1 < \ell_2 < \dots < \ell_J = \ell$ where J is the number of levels, and esti-
 325 mating a sequence of conditional probabilities $\mathbb{P}\{S(X) > \ell_{j+1} | S(X) > \ell_j\} =: p(\ell_{j+1} | \ell_j)$,
 326 which all have moderate magnitudes and are expected to be easier to estimate. Their
 327 product provides an estimate for the target probability:

$$\hat{p}_{\text{SS}}(\ell) = \hat{p}(\ell_1) \hat{p}(\ell_2 | \ell_1) \dots \hat{p}(\ell_J | \ell_{J-1}). \quad (5)$$

328 The first term can be estimated by vanilla Monte Carlo: generate N samples X_1, \dots, X_N ,
 329 and attach unit weights to each: $W_n = 1$ for $n = 1, \dots, N$. Rank the samples by S so
 330 that $S(X_{(1)}) \leq S(X_{(2)}) \leq \dots \leq S(X_{(N)})$, and let $\hat{p}(\ell_1) = (N - \kappa_1)/N$, where κ_1 is
 331 chosen so that $S(X_{(\kappa_1)}) \leq \ell_1 < S(X_{(\kappa_1+1)})$. The parameter κ_1 is the number of tra-
 332 jectories that are “killed” meaning they don’t appear in the first subset (see below). For
 333 the case of AMS, κ_1 is chosen as a parameter of the algorithm, and ℓ_1 is then set adap-
 334 tively as $\ell_1 = \frac{1}{2}[S(X_{(\kappa_1)}) + S(X_{(\kappa_1+1)})]$.

335 The second term $\hat{p}(\ell_2 | \ell_1)$ is estimated with a splitting strategy in which we focus
 336 in on the “subset” of samples that exceed the first threshold: $\{S(X) > \ell_1\}$ containing
 337 samples $X_{(i)}$ with $\kappa_1 < i \leq N$. To better sample this subset, we spawn additional sam-
 338 ples from it via a “Modified Metropolis algorithm”:

- 339 1. Initialize a list $\mathbb{X}_1 = \{X_{(\kappa_1+1)}, \dots, X_{(N)}\}$, which will eventually grow to a (user-
 340 chosen) size N_1 as well as a first-in-first-out queue \mathbb{Q} of the same elements but in
 341 a random order: the “parent queue”.
- 342 2. Pop \mathbb{Q} to yield the next parent X . Apply some small perturbation to X to gen-
 343 erate a new sample \tilde{X} , which itself is drawn from ρ but correlated to X . A gen-
 344 eral way to do this is with one step of the Metropolis-Hastings algorithm which
 345 involves an accept/reject step, but an easier approach is available in the partic-
 346 ular case of AMS as described in the next section.
- 347 3. Evaluate $S(\tilde{X})$. If it exceeds ℓ_1 , we have successfully generated a new sample from
 348 the subset. Accept the new sample, meaning insert \tilde{X} into both \mathbb{Q} and \mathbb{X}_1 and as-
 349 sign it a weight equal to that of its parent X . Otherwise, if $S(\tilde{X}) \leq \ell_1$, reject
 350 \tilde{X} . Re-insert X into \mathbb{Q} and add a copy of X to \mathbb{X}_1 . (In implementation, we don’t
 351 store two copies of the high-dimensional object X , but rather we assign a multi-
 352 plicity to each member and increment X ’s multiplicity by one.)
- 353 4. Repeat steps 2 and 3 until \mathbb{X}_1 has N_1 elements (counting multiplicity).
- 354 5. Multiply the weights of all members of \mathbb{X}_1 by a factor $(N - \kappa_1)/N_1$, which pre-
 355 serves the total weight N of the original ensemble $\{X_n\}_{n=1}^N$ while spreading that
 356 weight over more members.

357 Having expanded to N_1 samples from the subset $\{S(X) > \ell_1\}$, we can now pro-
 358 ceed to the next level and generate additional samples from the next subset $\{S(X) >$
 359 $\ell_2\}$ so that it contains N_2 samples, where ℓ_2 can be determined adaptively as an order
 360 statistic of \mathbb{X}_1 , i.e., the average of the κ_2 th and the (κ_2+1) th ranked values. The same
 361 procedure is repeated to generate the next subset \mathbb{X}_2 (and \mathbb{Q} is initialized with only unique
 362 elements, not counting multiplicity, in order to maintain as much diversity as possible).
 363 $\mathbb{X}_3, \mathbb{X}_4, \dots, \mathbb{X}_J$ are generated in the same fashion, until either a computational budget is
 364 reached, an ultimate target threshold is overcome, or some other halting criterion is met.

365 Ultimately we are left with a weighted ensemble $\{(X_1, W_1), \dots, (X_M, W_M)\}$, where
 366 $M = \kappa_1 + \kappa_2 + \dots + \kappa_J + N_J$. The sampling $\{S(X_m)\}_{m=1}^M$ is over-represented in the
 367 tails, but with correspondingly smaller weights there, and all weights sum to N . Any ex-

368 pectation of an observable $\Phi(x)$ can be estimated as

$$\mathbb{E}[\Phi(X)] = \int \Phi(x)\rho(x) dx \approx \hat{\Phi} = \frac{1}{N} \sum_{m=1}^M \Phi(X_m)W_m. \quad (6)$$

369 The SS algorithm will generally help to improve this estimate for functions Φ most sen-
 370 sitive to the tail region of $S(x)$, rather than its central bulk. In particular, setting $\Phi(x) =$
 371 $\mathbb{I}\{S(x) > \ell\}$, we recover the estimator $\hat{p}_{\text{SS}}(\ell)$:

$$\mathbb{E}[\mathbb{I}\{S(X) > \ell\}] = p(\ell) \approx \frac{1}{N} \sum_{m:S(X_m) > \ell} W_m = \hat{p}_{\text{SS}}(\ell). \quad (7)$$

372 An important set of algorithmic choices are the population parameters N, N_1, \dots, N_J ,
 373 the killing numbers $\kappa_1, \kappa_2, \dots, \kappa_J$, as well as the halting criterion which determines J . C erou
 374 et al. (2019) reviews theoretical bases for several different choices, but here for simplic-
 375 ity we opt for the same rule as used in Lestang et al. (2018): $\kappa_j = \kappa = 1$ (the ‘‘drop
 376 1’’ rule) and $N_j = N$ for all $j = 1, \dots, J$ (the population is replenished after each new
 377 level is set). Note that with $\kappa_j = 1$, only a single parent is selected from \mathbb{Q} at each round
 378 before the level is raised and the queue re-initialized.

379 3.1 Adaptive multilevel splitting (AMS)

380 AMS (in particular ‘‘trajectory AMS (TAMS)’’ in the nomenclature of Lestang et
 381 al. (2018)) can be seen as a special case of SS where each $X = \{X(t) : 0 \leq t \leq T\}$ is
 382 a length- T trajectory of a stochastic dynamical system, $S(X) = \max_{0 \leq t < T} R(X(t))$ for
 383 some time-dependent score function R , and with a particular choice for splitting trajec-
 384 tories. Trajectories are split by constructing a new forcing sequence $\tilde{\eta}(t)$ ($\tilde{\eta}_{1,2}(t)$ for our
 385 L96 model) to drive the child trajectory $\tilde{X}(t)$ starting from the old forcing sequence $\eta(t)$
 386 that drove the parent. First, copy the initial condition $\tilde{X}(0) = X(0)$. Then, copy $\tilde{\eta}(t) =$
 387 $\eta(t)$ up until some *split time* t_{sp} , which is chosen as first time $t_0(\ell)$ that the parent clears
 388 the threshold:

$$t_{\text{sp}} = t_0(\ell_1) = \min\{t \in [0, T] : R(X(t)) > \ell_1\}. \quad (8)$$

389 For following times $t \geq t_{\text{sp}}$, swap in a new and independent noise forcing sequence for
 390 $\tilde{\eta}(t)$. No Metropolis-style accept/reject step is needed for step (2) above; each newly sam-
 391 pled Brownian increment of $\tilde{\eta}(t)$ is drawn independently from $\mathcal{N}(0, \Delta t)$, and so $\tilde{\eta}(t)$ is
 392 a proper sample from the same noise-generating distribution as $\eta(t)$. Furthermore, the
 393 choice of $t_{\text{sp}} = t_0(\ell_1)$ guarantees $\tilde{X}(t) = X(t)$ for all $t \leq t_0(\ell_1)$, so that $S(\tilde{X}) > \ell_1$,
 394 and acceptance is guaranteed in step (3) as well.

395 The change in forcing for $t \geq t_{\text{sp}}$ will cause the child to diverge from the parent,
 396 producing a new—but correlated—sample (Fig. 1a). How correlated \tilde{X} is to its parent
 397 X depends on t_{sp} , with later t_{sp} implying a longer shared history and less independence.
 398 Applying the split at $t_{\text{sp}} = t_0(\ell)$ maximizes the independence of the child—and ulti-
 399 mately the diversity of the AMS ensemble—while guaranteeing $S(\tilde{X})$ exceeds ℓ_1 , and there-
 400 fore is accepted in the modified Metropolis Algorithm. The same procedure is carried
 401 out for every subsequent level.

402 We performed a sequence of AMS experiments with the following parameters:

- 403 1. Physical constants and timescales: $F_4 \in \{3, 1, 0.5, 0.25\}$ for the default case $a =$
 404 1 which gives the stochastically forced L96 model, and $F_4 = 3$ for the case $a =$
 405 0 which gives the OU process (shown in supplementary Figs. S1 and S2). We fix
 406 $F_0 = 6$, and $K = 40$ throughout, and set the time horizon to $T = 6$.
- 407 2. Ensemble sizes and population control: $N = N_j = 128$ and $\kappa_j = 1$ for $j =$
 408 1, 2, ..., $J = 896$ adhering to a fixed computational budget of 1024 time horizons

409 simulated. One additional halting criterion is imposed: if the population loses so
 410 much diversity that all active ensemble members descend from the same ances-
 411 tor, we terminate the algorithm early.

- 412 3. We repeat the whole procedure $M = 56$ times for each parameter set, with dif-
 413 ferent seeds for pseudo-random number generation. Each repetition will be called
 414 a “run” of AMS. Having multiple runs allows us to assess variance, and by using
 415 pooled estimates from all runs to hedge against stagnation within local optima of
 416 phase space in a particular run.

417 The initial N -member ensemble is generated as a sequence of consecutive blocks
 418 from a moderate initialization simulation of length $N \times T$ ($T = 6$ is the time horizon),
 419 after discarding the first 50 units as spinup. The spinup is initialized as $x_k(0) = F_0 +$
 420 $\frac{1}{1000} \sin\left(\frac{2\pi k}{K}\right)$. The random number generator used to create the noise forcing sequences
 421 $\eta_{1,2}(t)$ is seeded with $s \in \{0, \dots, M-1\}$, a different value for each AMS run with a fixed
 422 parameter set. The N initial blocks, although weakly correlated, comprise a sample from
 423 the steady-state distribution of the stochastic L96 system. Larger N reduces the vari-
 424 ability of the AMS results, but it also means more up-front cost and more rounds of split-
 425 ting needed to reach return times long enough to make the algorithm worthwhile.

426 We compare our results from AMS to a long DNS simulation of length 1.28×10^6
 427 (separate from the initialization), which is then further elongated by a factor of 40 (con-
 428 catenating all K timeseries end-to-end) into 5.12×10^7 , exploiting the statistical equiv-
 429 alence of all $K = 40$ sites of L96. This curve is our best estimate of ground truth. Note
 430 that the symmetry is only exploited to extend the DNS estimate, not the AMS estimate.
 431 In a climate model with zonal inhomogeneities, such as continents, it would be inappro-
 432 priate to aggregate different longitudes together.

433 Fig. 4a,b illustrates the effect of successive mutations over the course of the AMS
 434 algorithm, on the relatively easy test case with strong stochastic forcing, $F_4 = 3$ and
 435 $a = 1$ (the even easier case of $a = 0$ —the OU process with no interference from advection—
 436 is documented in Lestang et al. (2018) and included in supplementary Figs. S1 and S2
 437 for completeness). By design, the levels increase monotonically over the course of gen-
 438 erations and the descendant scores march upward, ultimately mutating the moderate an-
 439 cestor into an extreme descendant. Going beyond this successful “anecdote”, Fig. 5(a,b,c)
 440 confirm the benefit of AMS for a *statistically accurate* sampling of the distribution’s tails.
 441 Fig. 5a shows return period curves calculated with the modified block maximum method
 442 according to three datasets: the full weighted ensemble from AMS; the initialization (“Init”),
 443 consisting of N ensemble members per AMS run; and the long DNS simulation. The re-
 444 turn levels are interpolated onto a common logarithmically spaced grid of return peri-
 445 ods for easy comparison between the three data sources. Whereas return level estimates
 446 based on the initializations alone (blue) scatter considerably around the ground truth,
 447 AMS provides a tighter range of estimates (red) around the ground truth, and for ~ 3
 448 orders of magnitude-longer return periods, at only 8 times the cost of initialization (1024
 449 members from an initial 128). Moreover, each AMS run is ~ 5000 times less costly than
 450 the DNS run that gave the ground truth curve; altogether, the 56 AMS runs are ~ 100
 451 times less costly.

452 Another way of comparing AMS to DNS is by pooling together all members from
 453 the 56 ensembles and considering them as one larger ensemble of size $56 \times 1024 = 57344$.
 454 Fig. 5b shows the resulting statistics which have the advantage of extending to consid-
 455 erably longer return periods than the individual AMS runs. Here, as in Fig. 3, the er-
 456 ror bars are given by the basic bootstrap 95% confidence interval using 5000 bootstrap
 457 samples, but in the case of DNS (gray error bar), each bootstrap resampling contains
 458 only enough blocks to match the total simulation time used by AMS (including all in-
 459 dependent runs). This lets us compare the uncertainties fairly between the two meth-
 460 ods. In the case of AMS error bars, the members within a single run are not indepen-

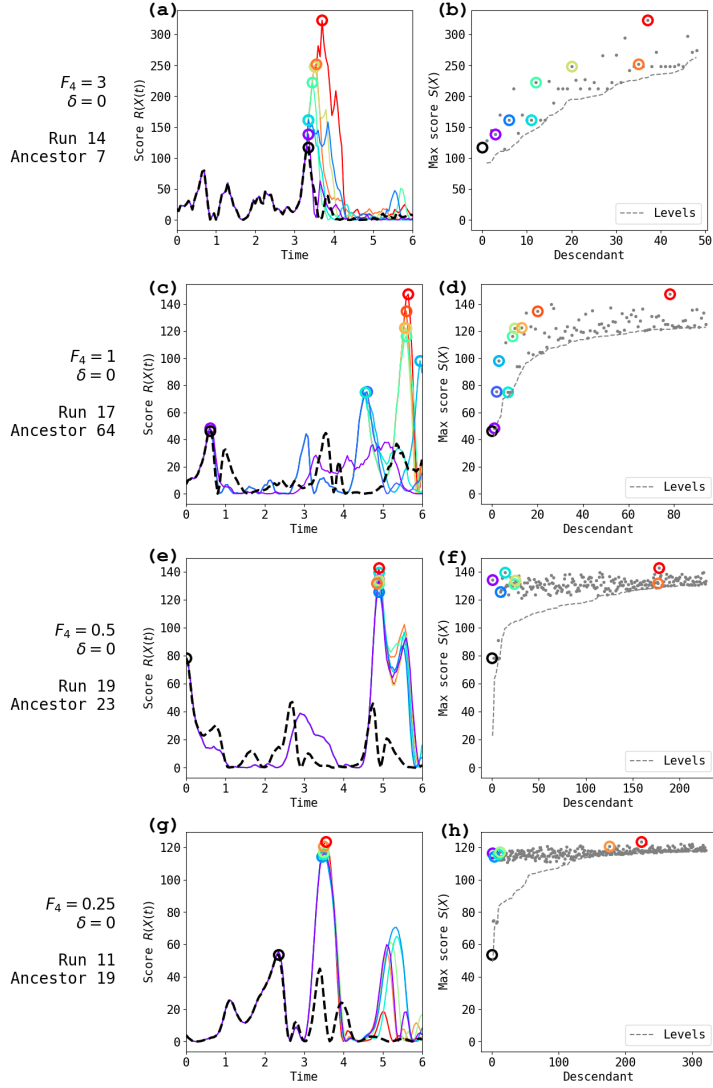


Figure 4. Scores for single ancestors and their descendants within the AMS algorithm (special case of TEAMS with $\delta = 0$). For each stochastic forcing amplitude, 56 independent runs of AMS were carried out (indexed 0-55) with $N = 128$ ensemble members (0-127). (a) Time-dependent score function $R(X(t))$ for the 7th initial ensemble member (ancestor) of run 14 for $F_4 = 3$. A black circle indicates the scalar score $S(X) = \max_t R(X(t))$. $R(X(t))$ and $S(X)$ are also shown for a single lineage (path down the family tree) in a sequence of brightening colors, ending with the highest scoring descendant’s score in red. (b) Scores in gray dots, with the horizontal axis numbering all descendants from ancestor 7 of run 14 for $F_4 = 3$. Colored circles indicate those descendants in the lineage from (a). The dashed gray curve indicates the levels ℓ from which each descendant was split. (c,e,g) are the same as (a), and (d,f,h) are the same as (b), but with stochastic forcing strength decreasing to $F_4 = 1, 0.5,$ and 0.25 respectively. In each case, the run and ancestor were hand-selected among the ancestors with the maximum boosting.

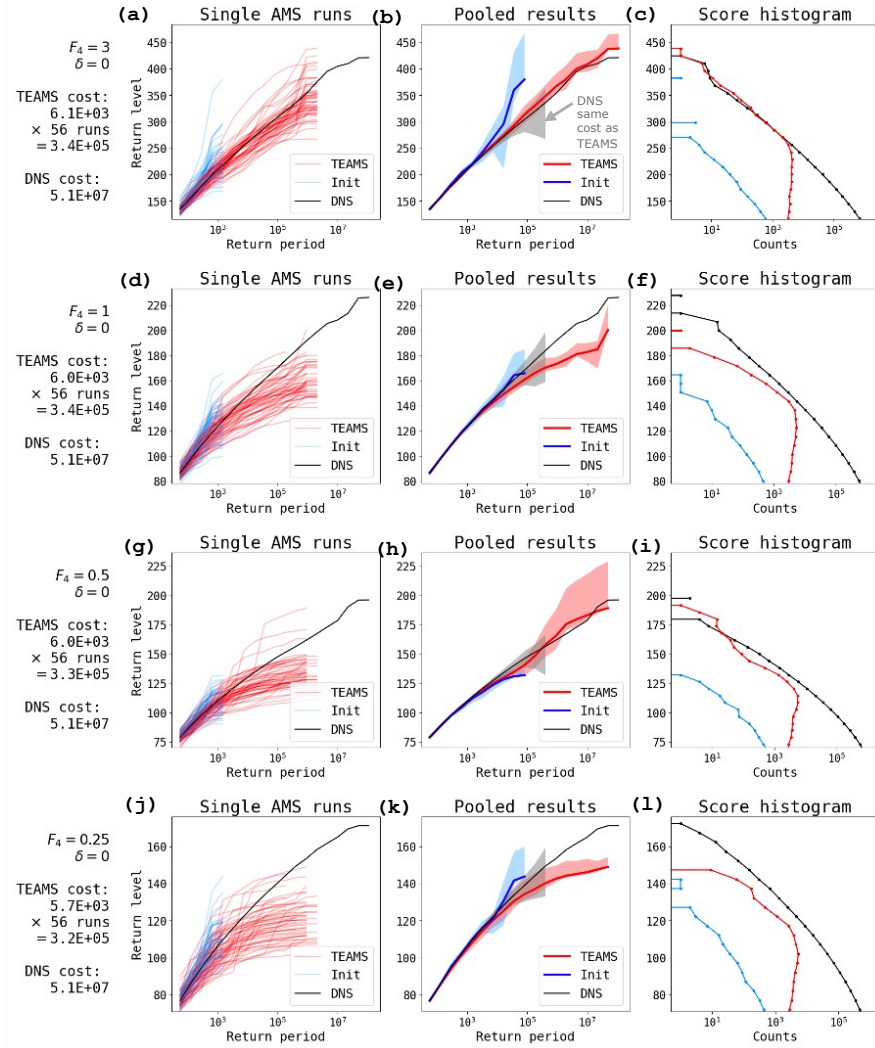


Figure 5. Performance of the AMS algorithm (special case of TEAMS with $\delta = 0$). (a) Return level vs. return period plots for $F_4 = 3$. Blue lines show estimates from the initial 128 members of each AMS run; red lines show estimates from the completed AMS runs; black line shows DNS. (b) Return level vs. return period for a pooled AMS ensemble containing all 56×1024 members. Blue and red envelopes indicate 95% confidence intervals (see text for details). Gray envelope is a 95% confidence interval based on subsets of DNS equal in total cost to the 56 AMS runs. Thus, the dashed red line and shading from AMS is of equal cost to the gray shading from DNS. (c) Unweighted histogram of scores for AMS initialization (blue), completed AMS (red), and DNS (black). Following rows are same as first row, but with noise decreasing to $F_4 = 1, 0.5$, and 0.25 , respectively. The slight variability in TEAMS costs listed to the left are due to the early halting criterion of one single ancestor remaining (see section 3).

461 dent of each other, and so we resample the AMS runs. That is, we sample the numbers
 462 $\{0, \dots, 55\}$ 5000 times with replacement, and for each resampling we pool together all mem-
 463 bers from the corresponding list of AMS runs, including repetitions. Fig. 5c shows the
 464 unweighted histogram of scores coming from the three data sources. The difference in
 465 shape of the AMS histogram compared to the DNS histogram demonstrates the main
 466 effect of AMS: to undersample the low end of the distribution and oversample the tail,
 467 shifting the computational burden to where it is more useful for sampling extremes.

468 We consider AMS to “win” over DNS if either of two criteria are met: (i) the AMS
 469 estimate remains close to the DNS (relative to error bar width) for return periods well
 470 beyond the AMS total simulation time T_{AMS} ; (ii) the AMS error bar is much smaller than
 471 the DNS error bar at T_{AMS} . Under strong stochastic forcing, AMS performs very well
 472 by both criteria, accurately (and confidently) estimating return periods as long as 10^7
 473 in the pooled estimate using only 3.4×10^5 time units of computation. This aligns with
 474 the demonstration in Lestang et al. (2018) for the OU process, and serves as a depart-
 475 ure point for our modification of the algorithm.

476 3.2 Failure of AMS in the regime of weak stochastic forcing

477 The story gets more complicated when the stochastic forcing is weak and nonlin-
 478 ear dynamics dominate. In deterministic chaos, perturbations grow exponentially with
 479 a rate inversely proportional to the *Lyapunov timescale*—at least, so long as the pertur-
 480 bations remain infinitesimal. Only after several elapsed Lyapunov times—what we call
 481 the *divergence timescale*, quantified further in section 4—do perturbations become large
 482 enough to be useful for splitting algorithms, but also at which size nonlinear effects take
 483 over. In contrast to deterministic chaos, white noise realizations diverge immediately.
 484 The stochastic L96 system inherits both behaviors to some extent, determined by the
 485 relative strength of stochastic forcing. Our main thesis is that when nonlinear dynam-
 486 ics dominate, and divergence time exceeds the duration of the event of interest, standard
 487 AMS is inadequate, but this can be remedied by adjusting the choice of splitting time
 488 t_{sp} as shown in the next section.

489 Fig. 4c-h show ancestors and descendants for AMS, analogous to Fig. 4a,b and with
 490 identical algorithmic parameters, but with decreasing levels of stochastic forcing: $F_4 =$
 491 $1, 0.5, 0.25$. For all four stochastic forcing strengths, ancestors can spawn more extreme
 492 descendants. However, there is a key difference between the strong- and weak-stochastic
 493 forcing regimes. With strong stochastic forcing $F_4 = 3$ (Fig. 4a,b), each descendant along
 494 the lineage improves upon the *same event*. In other words, the sequence of maximum
 495 scores comes from a peak in the timeseries for $R(X(t))$ that grows taller and taller, drift-
 496 ing only slightly forward in time. With weaker stochastic forcing (Fig. 4 c-d, e-f and es-
 497 pecially g-h), events tend to see only modest boosts from generation to generation. The
 498 only way for a child \tilde{X} to improve *substantially* over its parent X is by creating a whole
 499 new event—a new peak later in the time horizon—rather than building on an existing
 500 event. This happens because the stochastic forcing is too weak to open a large gap be-
 501 tween $R(\tilde{X}(t))$ and $R(X(t))$ during the short interval between the splitting time $t_0(\ell)$,
 502 when $R(X(t))$ first exceeds ℓ , and the peak $\text{argmax}_t R(X(t))$. The child ends up essen-
 503 tially replicating the parent’s peak, which is the same behavior illustrated schematically
 504 in Fig. 1a. The characteristic time scale of the peak (what we will call the event dura-
 505 tion) is set by the zonal propagation of waves, and this timescale is not long enough com-
 506 pared to the divergence time for AMS to work well. The same phenomenon was observed
 507 in Lestang et al. (2020): extreme spikes in the force on a body in a turbulent channel
 508 flow (see their Fig. 14) could not be boosted via AMS, which was attributed to the “sweep-
 509 ing” of vortices past the body. Similar reasoning holds for the zonal propagation of waves
 510 in L96 and the passage of midlatitude cyclones or fronts past a location in the midlat-
 511 itudes.

Fig. 5 summarizes the performance of AMS for different strengths of stochastic forcing. The suspicion of failure raised by Fig. 4 is confirmed by the clear degradation of performance as F_4 shrinks. In particular, the individual AMS return level curves tend to fall farther and farther underneath the true return level curves (left column of Fig. 5). There is a large scatter in the individual runs, and in the case $F_4 = 0.5$, a lucky few of the 56 runs salvage the pooled estimate for a decent approximation of the DNS return levels, but the width and asymmetry of the confidence intervals indicate the unreliability of this result (Fig. 5h). The problem becomes particularly acute as F_4 drops to 0.25, with the individual AMS runs barely improving upon the initial scores (Fig. 5j) and a large underestimate at longer return periods for the pooled estimate (Fig. 5k).

It thus appears that standard AMS is dead on arrival for cases where the divergence timescale is longer than the event duration. In principle, there is a canonical fix for this problem, namely to use a more intelligent score function than the quantity of interest $R(X(t))$ itself. The ideal such proxy is the *committor*: the probability, given an initial condition $X(t) = x$, that $R(X(s))$ will exceed ℓ at some time $s \in (t, T)$ before the time horizon ends. By definition, the committor incorporates information about the model state $X(t)$ that is not available from $R(X(t)) = x_0^2$, for example the speeds and magnitudes of different wave packets scattered across the domain that may all soon converge at $k = 0$ and result in an extreme burst of energy. The committor is an *optimal* score function for AMS in terms of minimizing the variance for $\hat{p}(\ell)$ (Lestang et al., 2018; Cérou et al., 2019; Lucente, Rolland, et al., 2022). Considerable research has recently pursued approximation strategies for the committor in various climate applications (e.g., Tantet et al., 2015; Finkel et al., 2021; Lucente, Herbert, & Bouchet, 2022; Miloshevich et al., 2023; Jacques-Dumas et al., 2023).

Unfortunately, these strategies all require either a high volume of training data—potentially canceling out the savings of a rare event algorithm, which is useful precisely in the low-data regime—or very specific knowledge of phase space geometry, such as a bistable structure, which is not typically available for realistic climate models. A second, related problem is that the optimality property only holds true for a single committor with a fixed threshold ℓ . What if we seek return periods for a whole range of thresholds? We would have to sacrifice the accuracy of some return periods in favor of others. Alternatively, we could use the committor for a single very high threshold ℓ_{\max} , but then even less training data would be available. Although it is interesting and worthwhile to search for committor functions based on traveling-wave dynamics, we leave that to future work, and in the next section we describe a simpler strategy to get around the stagnation issue seen in Fig. 4.

3.3 Trying-early adaptive multilevel splitting (TEAMS)

To address the failure of AMS in the nonlinear regime, we adjust $t_{\text{sp}} = t_\delta(\ell) =: t_0(\ell) - \delta$ by an *advance split time* $\delta > 0$, allowing some time for the child \tilde{X} to drift farther away from the parent and possibly achieve a higher maximum score. Indeed, ensemble boosting (Gessner et al., 2021) does exactly that, systematically applying perturbations every day from 19 to 7 days in advance of heat wave onset, although ensemble boosting does not by itself allow the calculation of return periods for the boosted events. When splitting early we lose the guarantee that $R(\tilde{X}(t))$ clears the current level ℓ (depicted schematically in Fig.1b), which is why we frame our modified algorithm using subset simulation (see section 3) which includes an accept/reject step: when a child fails to score higher than ℓ , it is discarded from the ensemble and its parent is duplicated instead (in other words, doubling its statistical weight). The resulting algorithm, which we call TEAMS (“trying-early adaptive multilevel splitting”), incurs additional cost due to rejected samples, but also gains back the ability to build significantly upon ancestral scores. One can interpret δ as setting the width of the proposal distribution, a key parameter in Markov chain Monte Carlo methods. A wider proposal allows the child to explore farther afield

564 from its parent, but increases the risk of rejection. Proposal width often has to be tuned
 565 carefully, and the sampling community has devoted substantial efforts to adaptively de-
 566 signing the proposal (Walter R. Gilks & Sahu, 1998; Andrieu & Thoms, 2008). Such meth-
 567 ods will surely prove useful for complex climate models, but in our present proof-of-concept
 568 study of the algorithm, we found approximately optimal δ values by exhaustive grid search
 569 for each noise level. Section 4 explains this procedure and shows that the optimal δ can
 570 be related to the error saturation timescale, a classical measure of predictability.

571 We performed a sequence of TEAMS experiments with $(F_4, \delta) \in \{3, 1, 0.5, 0.25\} \times$
 572 $\{0, 0.2, 0.4, \dots, 2.0\}$. We adjust the time horizon $T = 6 + \delta$ to give each parameter choice
 573 the same length of score to boost. All other parameters are as before for the AMS ex-
 574 periments.

575 Fig. 6 shows TEAMS in action for the same parameter sets from Fig. 4, but with
 576 (roughly optimal) advance splitting times $\delta = 0.0, 0.6, 1.0,$ and 1.4 for the decreasing
 577 noise levels (at $F_4 = 3$, $\delta = 0$ still works best, and panel (a) is the same as in Fig. 4a)).
 578 Note that the score functions $R(X(t))$ are only defined for times $t > \delta$, because if $t_0(\ell) <$
 579 δ then $t_\delta(\ell) < 0$, so we cannot apply the split early enough. This is implemented by
 580 setting the early scores to NaN, and lengthening the time horizon from T to $T + \delta$ as men-
 581 tioned above. We account for this extra cost in all the performance calculations to fol-
 582 low, but we omit the first δ time units from the plots. For all four stochastic forcing strengths,
 583 we see examples of children building significantly, and directly, upon a parent’s maxi-
 584 mum, without having to discover a new peak farther into the future. The values of the
 585 scores form continuous point clouds in panels (b,d,f,h), unlike the discrete horizontal bands
 586 appearing in Fig. 4(f,h) where $\delta = 0$ and stochastic forcing is weak. The negative side-
 587 effect is that many gray dots fall short of the gray dashed line, indicating a rejected sam-
 588 ple. Clearly, increasing δ brings both higher risk and higher reward.

589 Fig. 7 quantitatively confirms the hopeful suggestion of Fig. 6: that increasing δ
 590 can give TEAMS a speedup over DNS in the weak stochastic forcing regime. For all cases
 591 shown, TEAMS extends the estimated return period, *accurately*, well beyond the gray
 592 envelope which marks the limit achievable by an equal-cost run of DNS. The black ground
 593 truth curve remains within the 95% confidence band of TEAMS to return periods of \sim
 594 10^7 across all forcing levels. Simultaneously, the TEAMS confidence band is narrower
 595 than the DNS band.

596 Fig. 7 shows TEAMS gives a good estimate of the return values when all runs are
 597 pooled together, but that most individual TEAMS runs underestimate the true return
 598 values while a few overestimate them to allow for a good pooled estimate. As in Lucente,
 599 Rolland, et al. (2022), we can attribute this behavior to *apparent bias*, which is best ex-
 600 plained by analogy: an experiment consisting of 100 flips of a coin with $p = \mathbb{P}(\text{heads}) =$
 601 0.001 has a nine in ten chance of landing no heads, yielding a probability estimate $\hat{p} =$
 602 0 . But one experiment out of ten will yield $\hat{p} = 0.01$, a gross over-estimate, and only
 603 by pooling these two scenarios together can we see the estimator’s lack of bias. Unlike
 604 the coin-flipping experiment, TEAMS is designed to preferentially sample extreme val-
 605 ues, but a given AMS run for L96 may still get stuck in a local optimum yielding un-
 606 derestimated return values, especially if the stochastic forcing is too weak to jolt a tra-
 607 jectory out of it. Thus, pooling over multiple runs is especially crucial in the determi-
 608 nistic limit.

609 4 Optimizing advance split time

610 In this section, we explain how we determined optimal values of the advance split
 611 time δ using a simple exhaustive search. We then investigate the behavior of δ as a func-
 612 tion of stochastic forcing strength as a guide for choosing δ prior to running TEAMS on
 613 a more expensive model for which exhaustive search would not be feasible.

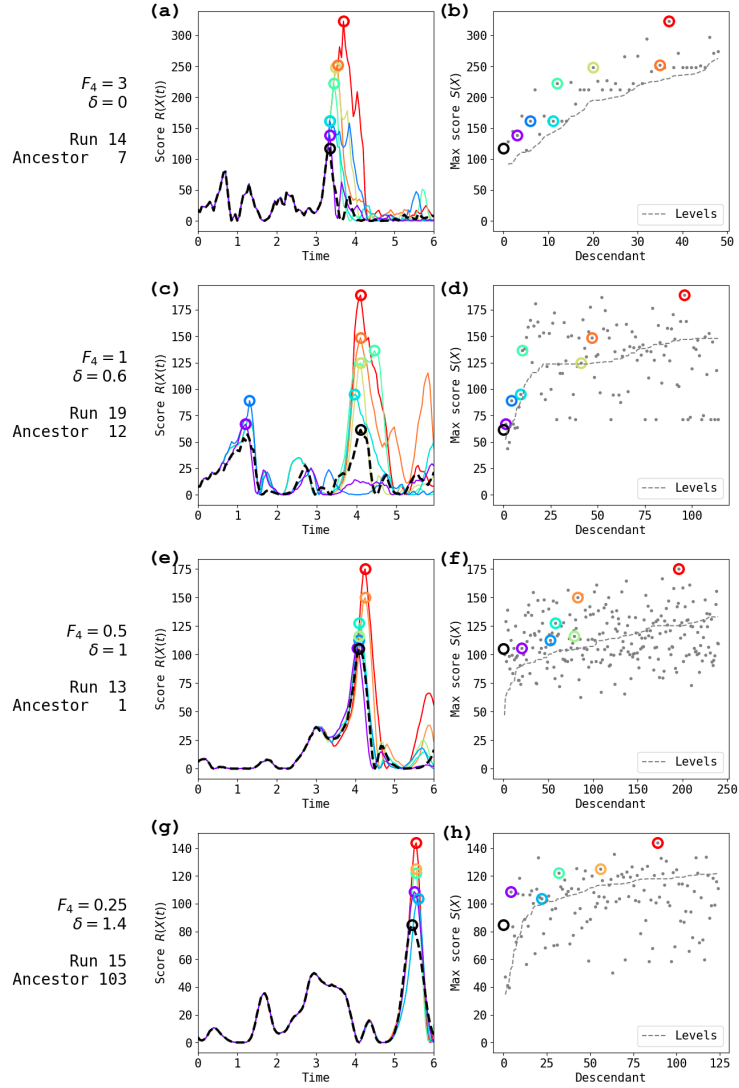


Figure 6. Scores for single ancestors and their descendants generated by the TEAMS algorithm: the same as Fig. 4 but with advance split times δ chosen to be approximately optimal for each noise level: $\delta = 0, 0.6, 1,$ and 1.4 for $F_4 = 3, 1, 0.5,$ and $0.25,$ respectively. Because $\delta = 0$ is optimal for $F_4 = 3,$ (a,b) is the same as Fig. 4a,b. Section 4 explains how the δ values were chosen.

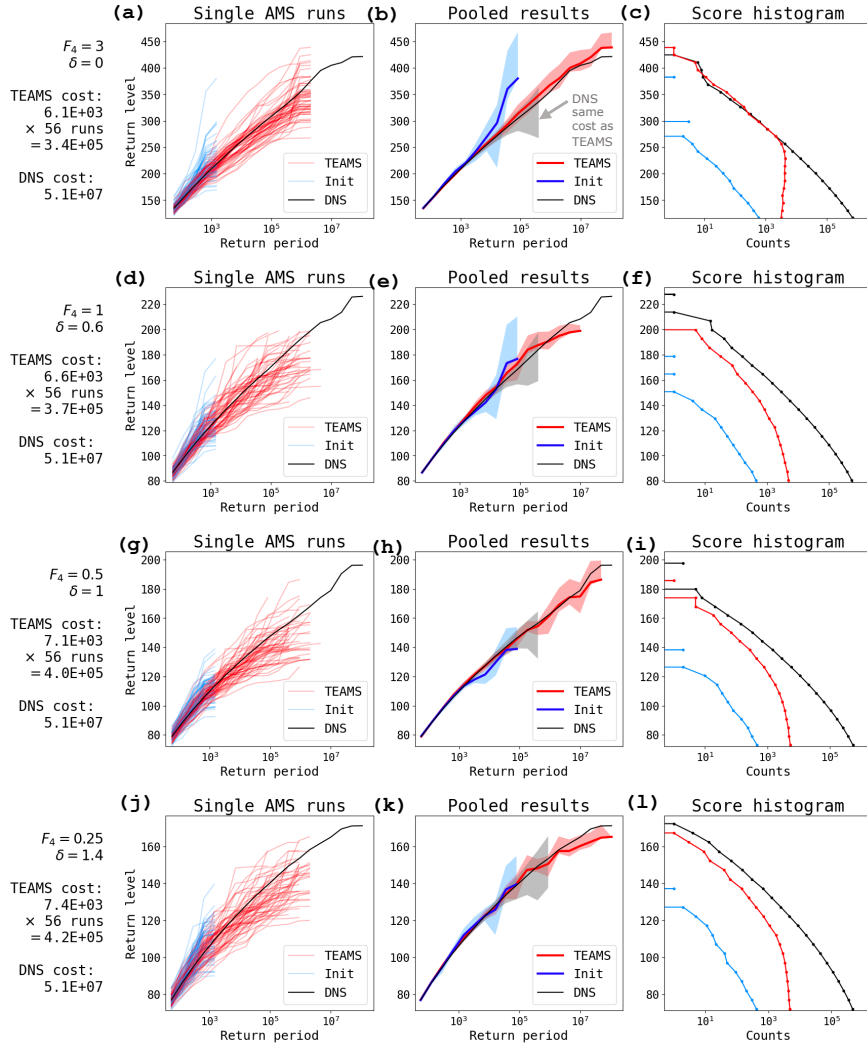


Figure 7. Performance of the TEAMS algorithm: the same as Fig. 5 but with advance split times δ chosen to be approximately optimal for each noise level: $\delta = 0, 0.6, 1, \text{ and } 1.4$ for $F_4 = 3, 1, 0.5, \text{ and } 0.25$, respectively. Because $\delta = 0$ is optimal for $F_4 = 3$, (a-c) are the same as Fig. 5a-c.

614

4.1 Exhaustive search

615

616

We selected the “optimal” δ values based on two simple performance metrics, which are plotted in Fig. 8.

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

1. Return level RMSE: the root-mean-square difference of return level between a TEAMS estimate (from a single run) and the DNS-determined ground truth, where the mean is taken over uniform bins in $\log \tau$ space. This metric is proportional to the L^2 -norm between a red line and the black line in the left columns of Figs. 5 and 7. In cases where the red line stops before the black line, it is extrapolated to longer return periods with a constant given by its maximum to penalize the algorithm getting stuck at a false upper bound. We calculate statistics of the return level RMSE across runs, including the mean and quantiles, which are displayed in Fig. 8(a,c,e,g). Note that these correspond to *percentile bootstrap* confidence intervals (Wasserman, 2004), as opposed to the *basic bootstrap* confidence intervals shown in Figs. 5 and 7. Here we use the percentile bootstrap as a means of sensitivity analysis, to show the range of results that might occur due to sampling variability. The basic bootstrap, by contrast, is intended to bracket the ground truth of some parameter value. The return level RMSE can also be calculated for the pooled estimate, and it shows similar but noisier trends.
2. Mean family gain: the maximum improvement (difference in scores) from ancestor to descendant over all N ancestors, averaged over the 56 runs. This does not measure statistical accuracy, but only the consistent ability to generate extreme events out of moderate events. Fig. 8 (b,d,f,h) shows mean family gain. Other metrics of gain, such as the maximum descendant score minus the maximum ancestral score (not necessarily from the same family tree) yield very similar trends with δ , albeit different absolute values.

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

A good choice of δ should have a small return level RMSE and a large mean family gain. Based on both performance metrics, we selected optimal $\delta = 0, 0.6, 1, 1.4$ for $F_4 = 3, 1, 0.5, 0.25$, respectively. These optimal values are marked with vertical gray lines in Fig. 8, and they are used in Figs. 6 and 7. For $F_4 = 0.5$, the two metrics gave slightly different optimal values ($\delta = 1.2$ for return level RMSE or $\delta = 1$ for mean family gain); we chose $\delta = 1$ because it gave the better pooled estimate. We emphasize that the optimal values are only discernible by averaging over many independent runs. For completeness, we display all 44 return level vs. return period plots (4 values of $F_4 \times 11$ values of δ) in the supplement. In general, shifting the optimal δ by ± 0.2 doesn’t change the results qualitatively, but larger shifts can degrade performance. The absolute values of errors should not be compared between stochastic forcing levels, since each has its own statistical steady state. Rather, the important takeaway is the increase in optimal δ as the stochastic forcing weakens. Indeed, in the singular limit of zero stochastic forcing the advance split time must necessarily go to infinity to have any effect at all, and initial condition perturbations would be needed to split trajectories.

654

655

656

657

658

To summarize, we have found that some choices of δ can make TEAMS effective where AMS is not effective, and that the optimal δ increases as stochastic forcing magnitude decreases. In the next section we relate this behavior to the predictability time, which points toward a cheap method to estimate an optimal—or at least, reasonably performant— δ , without having to repeatedly run TEAMS.

659

4.2 Relation between optimal advance time and error saturation timescales

660

661

662

663

Heuristically, we expect the optimal advance time δ to reflect the divergence timescale of perturbed trajectories that are introduced in splitting. Can this be related to classical predictability timescales? Lyapunov analysis describes perturbation growth by way of Lyapunov exponents and singular vectors (Cencini & Ginelli, 2013; Norwood et al.,

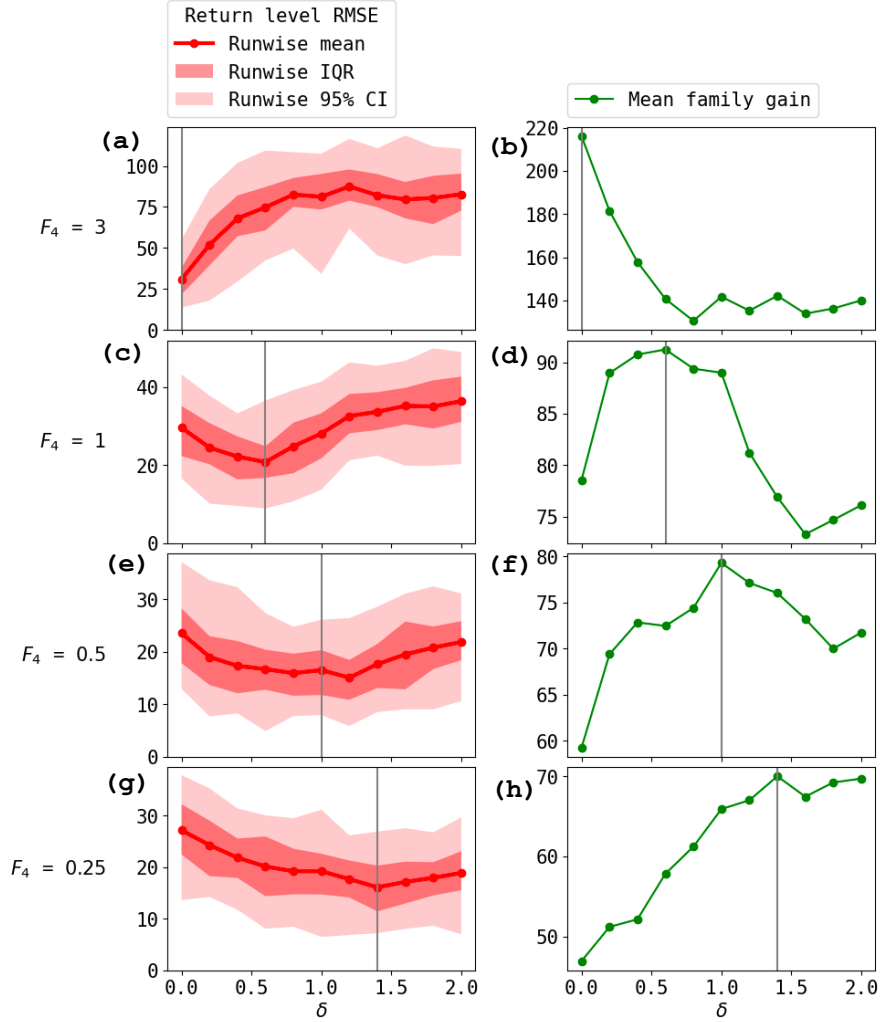


Figure 8. Performance of TEAMS as a function of advance split time δ and as measured by (a,c,e,g) return level RMSE and (b,d,f,h) mean family gain for $F_4 =$ (a,b) 3, (c,d) 1, (e,f) 0.5, and (g,h) 0.25. Return level RMSE is computed separately for each run. Thick red lines show the average over runs, and red envelopes show the quantile ranges 25%-75% (or interquartile range, IQR) and 2.5%-97.5% across the 56 runs. Mean family gain is maximum gain in score within a single family averaged over the 56 runs. Vertical gray lines show the optimal values of δ used in Figs. 6 and 7.

2013; Pazo et al., 2010; Maiocchi et al., 2024), but it applies to the regime of *infinitesimal* perturbations. The kind of perturbations we strive for in rare event sampling are finite and nonlinear, turning peaks into substantially larger peaks as in Figs. 4, 6. “Finite size Lyapunov exponents” (FSLEs) (Boffetta et al., 1998; Cencini & Vulpiani, 2013) are closer to what we need, generalizing the Lyapunov exponent to be dependent on an initial error amplitude. Typically, error grows in two stages: first exponentially, during which the FSLE equals the leading Lyapunov exponent, and then diffusively (scaling as a power law with time), during which the FSLE declines. The divergence timescale is bounded below by this change point, which approaches zero as stochastic forcing becomes dominant: indeed, the variance of pure Brownian motion grows linearly in t immediately.

On the other hand, the optimal δ is bounded above by the *error saturation timescale*, when perturbed ensemble members become independent and inhabit totally different regions of the attractor. By then, the root-mean-square error (RMSE) of the ensemble will equal the root-mean-square distance (RMSD) between two randomly chosen points on the attractor. In climate models, the saturation timescale is a convenient and effective measure of predictability (Sheshadri et al., 2021). Clearly, δ must be chosen shorter than the time to saturation, since a child trajectory ought to take advantage of pre-existing maxima produced by its parent. To investigate this relationship, the following experiments measure time in terms of fraction of saturation.

For each F_4 considered, we ran a moderate-length control simulation $x(t)$ for $0 \leq t \leq 1050$ (discarding the first 50 as spinup), and measured the RMSD for this simulation. At initialization times 50, 70, 90, ..., 990 (48 in total) we branched a 16-member ensemble with identical initial conditions $x(t)$ but independent stochastic forcing realizations (a convenient feature of stochastic forcing is that errors grow even from perfect initial conditions, removing dependence on initial perturbation amplitude). We integrated each member for 15 time units, calculated RMSE as a function of time (separately for each ensemble), and inverted to find the times t_ϵ at which the fraction of saturation $\epsilon = \text{RMSE}/\text{RMSD}$ reached a given value. In other words, $\text{RMSE}(t_\epsilon) = \epsilon \times \text{RMSD}$. Finally, we take the average across initializations to get a single value \bar{t}_ϵ for each of several ϵ values. The total cost of this experiment is 1.2×10^4 time units, roughly equal to 1.5 runs of AMS and much cheaper than the 56-run pooled estimate. Moreover, halving the number of initializations used yields qualitatively similar results.

Fig. 9 shows timeseries of $x_0(t)$ (both control and perturbed) and error growth for two such ensembles from the high and low stochastic forcing cases. The time axis is truncated to 10 days past initialization. The early linear growth of ϵ vs. \bar{t}_ϵ indicates a steady decline in relative growth rate, meaning that the perturbations begin to enter the diffusive (sub-exponential) growth regime quite early. This is thanks to stochastic forcing, which is visible in the top row as the emergence of red members from the shadow of the control trajectory. As expected, the error growth is much faster for the higher value of stochastic forcing.

If the optimal δ could be predicted from the error growth rates alone, the TEAMS algorithm could be calibrated simply and cheaply before being deployed. Fig. 10 shows the time $\bar{t}_{3/8}$ when RMSE reaches a fixed fraction of RMSD (3/8) as compared to the optimal δ values determined from Fig. 8, as a function of the strength of stochastic forcing. We include results from forcing at wavenumbers $m = 1, 4, 7, 10$. There is an encouraging similarity between the dependence of optimal δ and $\bar{t}_{3/8}$ on stochastic forcing strength, suggesting that the fractional saturation time might be useful to provide an estimate for δ .

Another interesting and less obvious feature is the dependence on wavenumber of error growth (albeit a weak dependence): medium-length wave forcing ($m = 4$ and $m = 7$) drives error to saturation faster than very short ($m = 10$) or long ($m = 1$) wave forcing, which informed our choice of $m = 4$ throughout the TEAMS experiments. How-

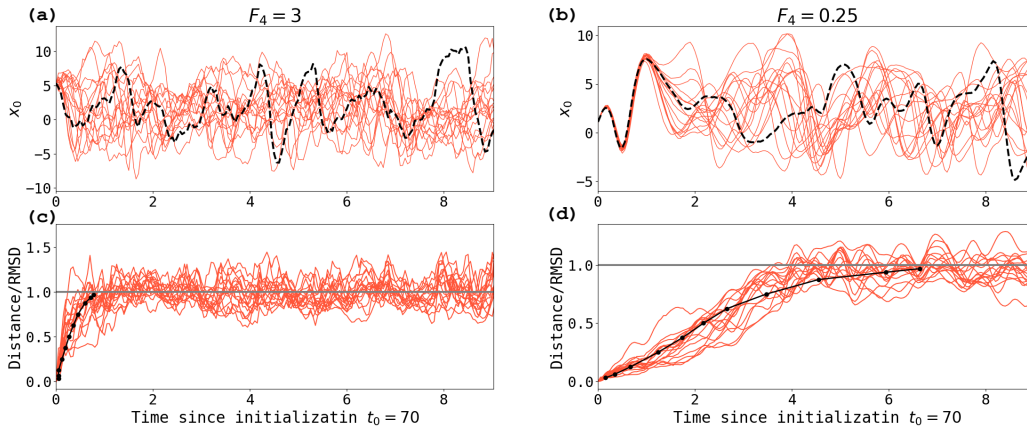


Figure 9. Growth of perturbations in the experiments described in subsection 4.2 for one representative initialization time $t_0 = 70$ and two values of the stochastic forcing: (a,c) $F_4 = 3$ and (b,d) $F_4 = 0.25$. (a,b) show $x_0(t)$ for the control simulation (black) and 16 simulations with the same initial condition but different white-noise forcing realizations (red). (c,d) show Euclidean distance between each ensemble member to the control as a fraction of RMSD versus time (red), and the fraction of saturation RMSE/RMSD versus the time until each ϵ value is reached averaged across all initializations and ensemble members (black), i.e., ϵ vs. \bar{t}_ϵ . Dots indicate $\epsilon = 1/32, 1/16, 1/8, 1/4, 3/8, 1/2$, and these same values reflected about $1/2$.

716 ever, the variability due to initial conditions (indicated by $\pm 1\sigma$ error bars) tend to exceed
 717 systematic differences between wavenumbers. This variability reflects a distribu-
 718 tion of divergence timescales across the attractor, which was also found to be quite het-
 719 erogeneous in Maiocchi et al. (2024) (there measured by Lyapunov exponents). It also
 720 suggests that the best strategy may be to not fix a single δ , but to allow the algorithm
 721 to adaptively set a δ , or sample from a range, to account for differing divergence timescales
 722 between different initial conditions, and this could be investigated in future work.

723 5 Conclusions and Outlook

724 A vexing challenge in climate science is reliably quantifying the probability of extreme
 725 weather events, which are fundamentally difficult to characterize because of data
 726 scarcity. Among various competing strategies, rare event algorithms hold several key ad-
 727 vantages, chiefly (i) access to dynamical samples of the events, rather than just return
 728 period curves which extreme value theory might provide, and (ii) more statistical rigor
 729 than storyline-based approaches like “ensemble boosting” (Gessner et al., 2021), thanks
 730 to careful re-weighting of cloned trajectories. Inspired by recent successes of rare event
 731 algorithms on long-lasting heat waves (Ragone et al., 2018) and idealized models of regime
 732 transitions (Lucente, Rolland, et al., 2022; Jacques-Dumas et al., 2023), we have inves-
 733 tigated the ability of a particular algorithm, adaptive multilevel splitting (AMS) to sam-
 734 ple extreme events of a different character: intermittent, short-lived bursts of energy in
 735 the Lorenz-96 model which have some similar characteristics as extreme daily rain or wind
 736 associated with midlatitude cyclones.

737 Even in this simple model, we have elucidated some key obstacles that hinder rare
 738 event splitting algorithms on sudden, short-lived events, and furthermore taken some steps

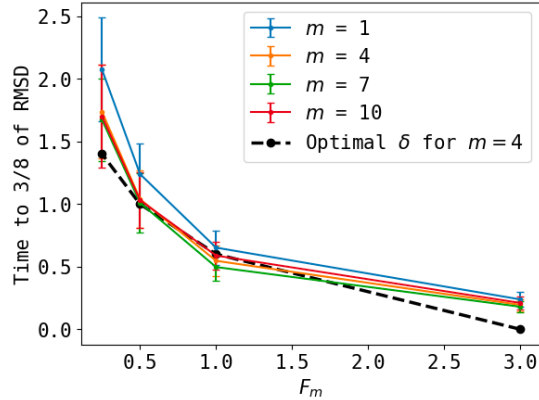


Figure 10. Time $\overline{t_{3/8}}$ until the perturbations described in subsection 4.2 reach a fixed fraction (3/8) of RMSD as a function of stochastic forcing strength F_m for different wavenumbers m . Error bars are ± 1 standard deviation of the distribution over different initial conditions. Optimized values of δ (determined from the performance metrics in Fig. 8) are shown in the black dashed line for $m = 4$.

739 to overcome them. AMS sets up a sequence of thresholds for an observable of interest
 740 and estimates conditional exceedance probabilities in stages by cloning and perturbing
 741 “successful” ensemble members when they cross a threshold, to generate new “success-
 742 ful” samples. This simple prescription suffers a fatal problem when the events are short-
 743 lived compared to the divergence timescale (how long it takes a perturbation to grow ap-
 744 preciablely): a perturbed ensemble member essentially replicates its parent’s success, and
 745 doesn’t develop its own history until after the event is over. Neither the magnitude nor
 746 the diversity of rare event samples is enhanced. To fix this problem, we have drawn in-
 747 spiration from ensemble boosting to apply a perturbation *in advance* of the rare event
 748 by some lead time δ . But we have also retained rigorous statistics for these “storylines”
 749 by exploiting a more general rare event algorithm, subset simulation (SS), of which AMS
 750 is only a special case. We name the resulting algorithm “trying-early AMS” (TEAMS)
 751 and demonstrate its success in sampling the tails of the rare event distribution more ef-
 752 ficiently than direct numerical simulation can do, despite an extra computational cost
 753 due to rejected samples.

754 Our study is a proof of concept that suggests a path forward, but with some open
 755 questions and directions for improvement, which we summarize here:

- 756 • The most crucial algorithmic parameter is the advance split time, δ , which is equiv-
 757 alent to a proposal distribution width. Our grid search for optimal δ , though not
 758 a scalable solution, demonstrates a relationship with the time over which pertur-
 759 bations grow to a fraction of saturation. An important goal for future work is to
 760 assess this result for other underlying models such as general circulation models
 761 or for other error growth metrics. Given the localized nature of our observable (x_0^2
 762 is the energy at a single longitude site), it is interesting that a *global* Euclidean
 763 metric correlates with the optimal δ . Weighting the metric more heavily for grid
 764 points near $k = 0$ might further improve this relationship.
- 765 • The weak stochastic forcing limit $F_m \rightarrow 0$ is important to confront for climate
 766 models, which may be more practical to perturb just at the splitting time rather
 767 than continuously at every time step, especially if the climate model is not already
 768 equipped with a stochastic subgrid parameterization. In the TEAMS framework,

769 this would translate to perturbing a simulation at a lead time δ ahead of the event,
 770 but not at all following times. Perturbing at just one time makes a given pertur-
 771 bation magnitude less powerful—but also opens up interesting possibilities such
 772 as the use of deterministic optimization strategies to more efficiently discover the
 773 most extreme event possible from a given initial condition. For example, some di-
 774 rections of perturbation (singular vectors) grow much faster than others, a fact
 775 which has informed ensemble design in operational weather forecasting (Palmer
 776 & Zanna, 2013), and could be used to further improve the algorithm. Methods
 777 such as conditional nonlinear optimal perturbation (Wang et al., 2020, and ref-
 778 erences therein), generalized stability theory (Farrell & Ioannou, 1996), and large
 779 deviation theory (Dematteis et al., 2018, 2019; Schorlepp et al., 2023) may prove
 780 useful for this task.

- 781 • Related to the previous point, it is desirable to have greater efficiency with sam-
 782 ples in order to deploy rare event algorithms at scale. For example, we should not
 783 simply discard rejected samples, but rather learn from their “mistakes” to design
 784 better perturbations. Frameworks like Bayesian optimization and adaptive impor-
 785 tance sampling based on model reduction have been developed for this task, and
 786 have been used in safety assessment for reliability engineering (e.g., Cousins & Sap-
 787 sisis, 2014; X. Huang et al., 2016; Mohamad & Sapsis, 2018; Sapsis, 2020; Uribe et
 788 al., 2021; Zhang et al., 2022).

789 Rare event algorithms represent a new way to allocate computational resources to
 790 where they matter most. To realize their considerable potential for efficiency gains, we
 791 have taken one of the necessary steps to make them flexible enough to target intermit-
 792 tent, localized, transient events that characterize phenomena such as heavy precipita-
 793 tion in complex global climate models. The Lorenz-96 model is an invaluable prototype
 794 as a cheap system that poses similar algorithmic challenges. Forthcoming papers will use
 795 the insight gained here as a stepping stone to more complex and realistic models.

796 Data availability statement

797 The software to simulate and sample extreme events in Lorenz-96 using TEAMS
 798 is available in a public Zenodo repository at <https://zenodo.org/doi/10.5281/zenodo.10608187>.
 799 Interested readers are encouraged to try out the algorithm on other systems of interest,
 800 and should not hesitate to contact J. F. for assistance.

801 Acknowledgments

802 This research is part of the MIT Climate Grand Challenge on Weather and Climate Ex-
 803 tremes. It received support by the generosity of Eric and Wendy Schmidt by recommen-
 804 dation of Schmidt Sciences as part of its Virtual Earth System Research Institute (VESRI).
 805 Computations for this research were carried out on the MIT Engaging cluster.

806 References

- 807 Abbot, D. S., Webber, R. J., Hadden, S., Seligman, D., & Weare, J. (2021, Dec).
 808 Rare Event Sampling Improves Mercury Instability Statistics. *The Astrophys-*
 809 *ical Journal*, *923*(2), 236. Retrieved from [https://dx.doi.org/10.3847/](https://dx.doi.org/10.3847/1538-4357/ac2fa8)
 810 [1538-4357/ac2fa8](https://dx.doi.org/10.3847/1538-4357/ac2fa8) doi: 10.3847/1538-4357/ac2fa8
- 811 Adachi, S. A., & Tomita, H. (2020). Methodology of the Constraint Condition
 812 in Dynamical Downscaling for Regional Climate Evaluation: A Review.
 813 *Journal of Geophysical Research: Atmospheres*, *125*(11), e2019JD032166.
 814 Retrieved from [https://agupubs.onlinelibrary.wiley.com/doi/abs/](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019JD032166)
 815 [10.1029/2019JD032166](https://doi.org/10.1029/2019JD032166) (e2019JD032166 10.1029/2019JD032166) doi:
 816 <https://doi.org/10.1029/2019JD032166>

- 817 Andrieu, C., & Thoms, J. (2008, Dec). A tutorial on adaptive MCMC. *Statistics*
818 *and Computing*, 18(4), 343-373. Retrieved from [https://doi.org/10.1007/](https://doi.org/10.1007/s11222-008-9110-y)
819 [s11222-008-9110-y](https://doi.org/10.1007/s11222-008-9110-y) doi: 10.1007/s11222-008-9110-y
- 820 Au, S.-K., & Beck, J. L. (2001). Estimation of small failure probabilities in
821 high dimensions by subset simulation. *Probabilistic Engineering Mechan-*
822 *ics*, 16(4), 263-277. Retrieved from [https://www.sciencedirect.com/](https://www.sciencedirect.com/science/article/pii/S0266892001000194)
823 [science/article/pii/S0266892001000194](https://www.sciencedirect.com/science/article/pii/S0266892001000194) doi: [https://doi.org/10.1016/](https://doi.org/10.1016/S0266-8920(01)00019-4)
824 [S0266-8920\(01\)00019-4](https://doi.org/10.1016/S0266-8920(01)00019-4)
- 825 Baars, S., Castellana, D., Wubs, F., & Dijkstra, H. (2021). Application of adap-
826 tive multilevel splitting to high-dimensional dynamical systems. *Jour-*
827 *nal of Computational Physics*, 424, 109876. Retrieved from [https://](https://www.sciencedirect.com/science/article/pii/S0021999120306501)
828 www.sciencedirect.com/science/article/pii/S0021999120306501 doi:
829 <https://doi.org/10.1016/j.jcp.2020.109876>
- 830 Boffetta, G., Giuliani, P., Paladin, G., & Vulpiani, A. (1998). An Exten-
831 sion of the Lyapunov Analysis for the Predictability Problem. *Jour-*
832 *nal of the Atmospheric Sciences*, 55(23), 3409 - 3416. Retrieved from
833 [https://journals.ametsoc.org/view/journals/atsc/55/23/1520-0469](https://journals.ametsoc.org/view/journals/atsc/55/23/1520-0469_1998_055_3409_aeotla_2.0.co_2.xml)
834 [_1998_055_3409_aeotla_2.0.co_2.xml](https://journals.ametsoc.org/view/journals/atsc/55/23/1520-0469_1998_055_3409_aeotla_2.0.co_2.xml) doi: [https://doi.org/10.1175/](https://doi.org/10.1175/1520-0469(1998)055(3409:AEOTLA)2.0.CO;2)
835 [1520-0469\(1998\)055\(3409:AEOTLA\)2.0.CO;2](https://doi.org/10.1175/1520-0469(1998)055(3409:AEOTLA)2.0.CO;2)
- 836 Bouchet, F., Rolland, J., & Simonnet, E. (2019, Feb). Rare Event Algorithm
837 Links Transitions in Turbulent Flows with Activated Nucleations. *Phys. Rev.*
838 *Lett.*, 122, 074502. Retrieved from [https://link.aps.org/doi/10.1103/](https://link.aps.org/doi/10.1103/PhysRevLett.122.074502)
839 [PhysRevLett.122.074502](https://link.aps.org/doi/10.1103/PhysRevLett.122.074502) doi: 10.1103/PhysRevLett.122.074502
- 840 Bucklew, J. A. (2004). *Introduction to Rare Event Simulation* (1st ed.). Springer
841 New York, NY. doi: <https://doi.org/10.1007/978-1-4757-4078-3>
- 842 Buizza, R., Milleer, M., & Palmer, T. N. (1999). Stochastic representation of
843 model uncertainties in the ECMWF ensemble prediction system. *Quar-*
844 *terly Journal of the Royal Meteorological Society*, 125(560), 2887-2908. Re-
845 trieved from [https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.49712556006)
846 [qj.49712556006](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.49712556006) doi: <https://doi.org/10.1002/qj.49712556006>
- 847 Cencini, M., & Ginelli, F. (2013, Jun). Lyapunov analysis: from dynamical systems
848 theory to applications. *Journal of Physics A: Mathematical and Theoretical*,
849 46(25), 250301. Retrieved from [https://dx.doi.org/10.1088/1751-8113/](https://dx.doi.org/10.1088/1751-8113/46/25/250301)
850 [46/25/250301](https://dx.doi.org/10.1088/1751-8113/46/25/250301) doi: 10.1088/1751-8113/46/25/250301
- 851 Cencini, M., & Vulpiani, A. (2013, Jun). Finite size Lyapunov exponent: review
852 on applications. *Journal of Physics A: Mathematical and Theoretical*, 46(25),
853 254019. Retrieved from [https://dx.doi.org/10.1088/1751-8113/46/25/](https://dx.doi.org/10.1088/1751-8113/46/25/254019)
854 [254019](https://dx.doi.org/10.1088/1751-8113/46/25/254019) doi: 10.1088/1751-8113/46/25/254019
- 855 Coles, S. (2001). *An introduction to statistical modeling of extreme values* (1st ed.).
856 Springer. doi: 10.1007/978-1-4471-3675-0
- 857 Cousins, W., & Sapsis, T. P. (2014). Quantification and prediction of ex-
858 treme events in a one-dimensional nonlinear dispersive wave model. *Phys-*
859 *ica D: Nonlinear Phenomena*, 280-281, 48-58. Retrieved from [https://](https://www.sciencedirect.com/science/article/pii/S016727891400092X)
860 www.sciencedirect.com/science/article/pii/S016727891400092X doi:
861 <https://doi.org/10.1016/j.physd.2014.04.012>
- 862 Cérou, F., & Guyader, A. (2007). Adaptive Multilevel Splitting for Rare
863 Event Analysis. *Stochastic Analysis and Applications*, 25(2), 417-443.
864 Retrieved from <https://doi.org/10.1080/07362990601139628> doi:
865 [10.1080/07362990601139628](https://doi.org/10.1080/07362990601139628)
- 866 Cérou, F., Guyader, A., & Rousset, M. (2019). Adaptive multilevel splitting: His-
867 torical perspective and recent results. *Chaos: An Interdisciplinary Journal of*
868 *Nonlinear Science*, 29(4), 043108. Retrieved from [https://doi.org/10.1063/](https://doi.org/10.1063/1.5082247)
869 [1.5082247](https://doi.org/10.1063/1.5082247) doi: 10.1063/1.5082247
- 870 Dematteis, G., Grafke, T., & Vanden-Eijnden, E. (2018). Rogue waves and large
871 deviations in deep sea. *Proceedings of the National Academy of Sciences*,

- 872 115(5), 855-860. Retrieved from [https://www.pnas.org/doi/abs/10.1073/](https://www.pnas.org/doi/abs/10.1073/pnas.1710670115)
873 [pnas.1710670115](https://www.pnas.org/doi/abs/10.1073/pnas.1710670115) doi: 10.1073/pnas.1710670115
- 874 Dematteis, G., Grafke, T., & Vanden-Eijnden, E. (2019). Extreme Event Quantifi-
875 cation in Dynamical Systems with Random Components. *SIAM/ASA Journal*
876 *on Uncertainty Quantification*, 7(3), 1029-1059. Retrieved from [https://doi](https://doi.org/10.1137/18M1211003)
877 [.org/10.1137/18M1211003](https://doi.org/10.1137/18M1211003) doi: 10.1137/18M1211003
- 878 Dwyer, J. G., & O’Gorman, P. A. (2017). Changing duration and spatial ext-
879 tent of midlatitude precipitation extremes across different climates. *Geo-*
880 *physical Research Letters*, 44(11), 5863-5871. Retrieved from [https://](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2017GL072855)
881 agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2017GL072855 doi:
882 <https://doi.org/10.1002/2017GL072855>
- 883 Emanuel, K. (2021). Response of Global Tropical Cyclone Activity to Increasing
884 CO₂: Results from Downscaling CMIP6 Models. *Journal of Climate*, 34(1), 57
885 - 70. Retrieved from [https://journals.ametsoc.org/view/journals/clim/](https://journals.ametsoc.org/view/journals/clim/34/1/jcliD200367.xml)
886 [34/1/jcliD200367.xml](https://journals.ametsoc.org/view/journals/clim/34/1/jcliD200367.xml) doi: <https://doi.org/10.1175/JCLI-D-20-0367.1>
- 887 Farrell, B. F., & Ioannou, P. J. (1996). Generalized Stability Theory. Part I:
888 Autonomous Operators. *Journal of Atmospheric Sciences*, 53(14), 2025 -
889 2040. Retrieved from [https://journals.ametsoc.org/view/journals/](https://journals.ametsoc.org/view/journals/atsc/53/14/1520-0469_1996_053_2025_gstpia_2_0_co_2.xml)
890 [atsc/53/14/1520-0469_1996_053_2025_gstpia_2_0_co_2.xml](https://journals.ametsoc.org/view/journals/atsc/53/14/1520-0469_1996_053_2025_gstpia_2_0_co_2.xml) doi:
891 [10.1175/1520-0469\(1996\)053<2025:GSTPIA>2.0.CO;2](https://doi.org/10.1175/1520-0469(1996)053<2025:GSTPIA>2.0.CO;2)
- 892 Finkel, J., Webber, R. J., Gerber, E. P., Abbot, D. S., & Weare, J. (2021). Learn-
893 ing Forecasts of Rare Stratospheric Transitions from Short Simulations.
894 *Monthly Weather Review*, 149(11), 3647 - 3669. Retrieved from [https://](https://journals.ametsoc.org/view/journals/mwre/149/11/MWR-D-21-0024.1.xml)
895 journals.ametsoc.org/view/journals/mwre/149/11/MWR-D-21-0024.1.xml
896 doi: 10.1175/MWR-D-21-0024.1
- 897 Gagne II, D. J., Christensen, H. M., Subramanian, A. C., & Monahan, A. H.
898 (2020). Machine Learning for Stochastic Parameterization: Generative Ad-
899 versarial Networks in the Lorenz '96 Model. *Journal of Advances in Mod-*
900 *eling Earth Systems*, 12(3), e2019MS001896. Retrieved from [https://](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019MS001896)
901 agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019MS001896
902 [\(e2019MS001896 10.1029/2019MS001896\)](https://doi.org/10.1029/2019MS001896) doi: [https://doi.org/10.1029/](https://doi.org/10.1029/2019MS001896)
903 [2019MS001896](https://doi.org/10.1029/2019MS001896)
- 904 Gálfi, V. M., Bódai, T., & Lucarini, V. (2017, Sep 06). Convergence of Extreme
905 Value Statistics in a Two-Layer Quasi-Geostrophic Atmospheric Model. *Com-*
906 *plexity*, 2017, 5340858. Retrieved from [https://doi.org/10.1155/2017/](https://doi.org/10.1155/2017/5340858)
907 [5340858](https://doi.org/10.1155/2017/5340858) doi: 10.1155/2017/5340858
- 908 Gálfi, V. M., Lucarini, V., Ragone, F., & Wouters, J. (2021, Jun 01). Applications
909 of large deviation theory in geophysical fluid dynamics and climate science. *La*
910 *Rivista del Nuovo Cimento*, 44(6), 291-363. Retrieved from [https://doi.org/](https://doi.org/10.1007/s40766-021-00020-z)
911 [10.1007/s40766-021-00020-z](https://doi.org/10.1007/s40766-021-00020-z) doi: 10.1007/s40766-021-00020-z
- 912 Gessner, C. (2022). *Physical storylines for very rare climate extremes* (Unpublished
913 doctoral dissertation). ETH Zurich.
- 914 Gessner, C., Fischer, E. M., Beyerle, U., & Knutti, R. (2021). Very Rare Heat
915 Extremes: Quantifying and Understanding Using Ensemble Reinitializa-
916 tion. *Journal of Climate*, 34(16), 6619 - 6634. Retrieved from [https://](https://journals.ametsoc.org/view/journals/clim/34/16/JCLI-D-20-0916.1.xml)
917 journals.ametsoc.org/view/journals/clim/34/16/JCLI-D-20-0916.1.xml
918 doi: 10.1175/JCLI-D-20-0916.1
- 919 Giardinà, C., Kurchan, J., & Peliti, L. (2006, Mar). Direct Evaluation of
920 Large-Deviation Functions. *Phys. Rev. Lett.*, 96, 120603. Retrieved from
921 <https://link.aps.org/doi/10.1103/PhysRevLett.96.120603> doi:
922 [10.1103/PhysRevLett.96.120603](https://doi.org/10.1103/PhysRevLett.96.120603)
- 923 Hu, G., Bódai, T., & Lucarini, V. (2019). Effects of stochastic parametrization on
924 extreme value statistics. *Chaos: An Interdisciplinary Journal of Nonlinear Sci-*
925 *ence*, 29(8), 083102. Retrieved from <https://doi.org/10.1063/1.5095756>
926 doi: 10.1063/1.5095756

- 927 Huang, W. K., Stein, M. L., McInerney, D. J., Sun, S., & Moyer, E. J. (2016).
 928 Estimating changes in temperature extremes from millennial-scale climate
 929 simulations using generalized extreme value (GEV) distributions. *Advances*
 930 *in Statistical Climatology, Meteorology and Oceanography*, 2(1), 79–103. Re-
 931 trieved from <https://ascmo.copernicus.org/articles/2/79/2016/> doi:
 932 10.5194/ascmo-2-79-2016
- 933 Huang, X., Chen, J., & Zhu, H. (2016). Assessing small failure probabili-
 934 ties by AK–SS: An active learning method combining Kriging and Sub-
 935 set Simulation. *Structural Safety*, 59, 86–95. Retrieved from [https://](https://www.sciencedirect.com/science/article/pii/S0167473016000035)
 936 www.sciencedirect.com/science/article/pii/S0167473016000035 doi:
 937 <https://doi.org/10.1016/j.strusafe.2015.12.003>
- 938 Huang, X., Swain, D. L., & Hall, A. D. (2020). Future precipitation increase
 939 from very high resolution ensemble downscaling of extreme atmospheric
 940 river storms in California. *Science Advances*, 6(29), eaba1323. Retrieved
 941 from <https://www.science.org/doi/abs/10.1126/sciadv.aba1323> doi:
 942 10.1126/sciadv.aba1323
- 943 Jacques-Dumas, V., van Westen, R. M., Bouchet, F., & Dijkstra, H. A. (2023).
 944 Data-driven methods to estimate the committor function in conceptual
 945 ocean models. *Nonlinear Processes in Geophysics*, 30(2), 195–216. Re-
 946 trieved from <https://npg.copernicus.org/articles/30/195/2023/> doi:
 947 10.5194/npg-30-195-2023
- 948 Kahn, H., & Harris, T. E. (1951). Estimation of particle transmission by random
 949 sampling. *National Bureau of Standards applied mathematics series*, 12, 27–
 950 30.
- 951 Kharin, V. V., Zwiers, F. W., Zhang, X., & Hegerl, G. C. (2007). Changes in
 952 Temperature and Precipitation Extremes in the IPCC Ensemble of Global
 953 Coupled Model Simulations. *Journal of Climate*, 20(8), 1419 - 1444. Re-
 954 trieved from [https://journals.ametsoc.org/view/journals/clim/20/8/](https://journals.ametsoc.org/view/journals/clim/20/8/jcli4066.1.xml)
 955 [jcli4066.1.xml](https://journals.ametsoc.org/view/journals/clim/20/8/jcli4066.1.xml) doi: <https://doi.org/10.1175/JCLI4066.1>
- 956 Krouma, M., Yiou, P., Déandreis, C., & Thao, S. (2022). Assessment of stochas-
 957 tic weather forecast of precipitation near European cities, based on analogs
 958 of circulation. *Geoscientific Model Development*, 15(12), 4941–4958. Re-
 959 trieved from <https://gmd.copernicus.org/articles/15/4941/2022/> doi:
 960 10.5194/gmd-15-4941-2022
- 961 Kästner, J. (2011). Umbrella sampling. *WIREs Computational Molecular Science*,
 962 1(6), 932–942. Retrieved from [https://wires.onlinelibrary.wiley.com/](https://wires.onlinelibrary.wiley.com/doi/abs/10.1002/wcms.66)
 963 [doi/abs/10.1002/wcms.66](https://wires.onlinelibrary.wiley.com/doi/abs/10.1002/wcms.66) doi: <https://doi.org/10.1002/wcms.66>
- 964 Lee, C.-Y., CaMargo, S. J., Sobel, A. H., & Tippet, M. K. (2020). Statisti-
 965 cal–Dynamical Downscaling Projections of Tropical Cyclone Activity in a
 966 Warming Climate: Two Diverging Genesis Scenarios. *Journal of Climate*,
 967 33(11), 4815 - 4834. Retrieved from [https://journals.ametsoc.org/](https://journals.ametsoc.org/view/journals/clim/33/11/jcli-d-19-0452.1.xml)
 968 [view/journals/clim/33/11/jcli-d-19-0452.1.xml](https://journals.ametsoc.org/view/journals/clim/33/11/jcli-d-19-0452.1.xml) doi: 10.1175/
 969 JCLI-D-19-0452.1
- 970 Lestang, T., Bouchet, F., & Lévêque, E. (2020). Numerical study of extreme
 971 mechanical force exerted by a turbulent flow on a bluff body by direct and
 972 rare-event sampling techniques. *Journal of Fluid Mechanics*, 895, A19. doi:
 973 10.1017/jfm.2020.293
- 974 Lestang, T., Ragone, F., Bréhier, C.-E., Herbert, C., & Bouchet, F. (2018, Apr).
 975 Computing return times or return periods with rare event algorithms. *Jour-
 976 nal of Statistical Mechanics: Theory and Experiment*, 2018(4), 043213.
 977 Retrieved from <https://doi.org/10.1088/1742-5468/aab856> doi:
 978 10.1088/1742-5468/aab856
- 979 Lorenz, E. N. (1996). Predictability: A problem partly solved. In *Proc. Seminar on*
 980 *predictability* (Vol. 1). Retrieved from [https://www.cambridge.org/core/](https://www.cambridge.org/core/books/abs/predictability-of-weather-and-climate/predictability-a)
 981 [books/abs/predictability-of-weather-and-climate/predictability-a](https://www.cambridge.org/core/books/abs/predictability-of-weather-and-climate/predictability-a)

- 982 ~~-problem-partly-solved/3221BDE379DEB669BA52C66263AF3206~~
- 983 Lorenz, E. N., & Emanuel, K. A. (1998). Optimal Sites for Supplementary Weather
984 Observations: Simulation with a Small Model. *Journal of the Atmospheric*
985 *Sciences*, 55(3), 399 - 414. Retrieved from [https://journals.ametsoc.org/
986 view/journals/atsc/55/3/1520-0469.1998.055.0399_osfsw0.2.0.co.2.xml](https://journals.ametsoc.org/view/journals/atsc/55/3/1520-0469.1998.055.0399_osfsw0.2.0.co.2.xml)
987 doi: 10.1175/1520-0469(1998)055<0399:OSFSWO>2.0.CO;2
- 988 Lucarini, V., Faranda, D., de Freitas, J. M. M., Holland, M., Kuna, T., Nicol, M.,
989 ... others (2016). *Extremes and recurrence in dynamical systems*. John Wiley
990 & Sons.
- 991 Lucente, D., Herbert, C., & Bouchet, F. (2022). Committed Functions for Cli-
992 mate Phenomena at the Predictability Margin: The Example of El Niño
993 Southern Oscillation in the Jin and Timmermann Model. *Journal of the*
994 *Atmospheric Sciences*. Retrieved from [https://journals.ametsoc.org/
995 view/journals/atsc/aop/JAS-D-22-0038.1/JAS-D-22-0038.1.xml](https://journals.ametsoc.org/view/journals/atsc/aop/JAS-D-22-0038.1/JAS-D-22-0038.1.xml) doi:
996 10.1175/JAS-D-22-0038.1
- 997 Lucente, D., Rolland, J., Herbert, C., & Bouchet, F. (2022, Aug). Coupling rare
998 event algorithms with data-based learned committed functions using the ana-
999 logue Markov chain. *Journal of Statistical Mechanics: Theory and Experiment*,
1000 2022(8), 083201. Retrieved from [https://dx.doi.org/10.1088/1742-5468/
1001 ac7aa7](https://dx.doi.org/10.1088/1742-5468/ac7aa7) doi: 10.1088/1742-5468/ac7aa7
- 1002 Maiocchi, C. C., Lucarini, V., Gritsun, A., & Sato, Y. (2024). Heterogeneity of the
1003 attractor of the Lorenz '96 model: Lyapunov analysis, unstable periodic orbits,
1004 and shadowing properties. *Physica D: Nonlinear Phenomena*, 457, 133970.
1005 Retrieved from [https://www.sciencedirect.com/science/article/pii/
1006 S016727892300324X](https://www.sciencedirect.com/science/article/pii/S016727892300324X) doi: <https://doi.org/10.1016/j.physd.2023.133970>
- 1007 Miloshevich, G., Cozian, B., Abry, P., Borgnat, P., & Bouchet, F. (2023, Apr). *Prob-*
1008 *abilistic forecasts of extreme heatwaves using convolutional neural networks*
1009 *in a regime of lack of data* (Vol. 8). American Physical Society. Retrieved
1010 from <https://link.aps.org/doi/10.1103/PhysRevFluids.8.040501> doi:
1011 10.1103/PhysRevFluids.8.040501
- 1012 Mohamad, M. A., & Sapsis, T. P. (2018). Sequential sampling strategy for
1013 extreme event statistics in nonlinear dynamical systems. *Proceedings*
1014 *of the National Academy of Sciences*, 115(44), 11138-11143. Retrieved
1015 from <https://www.pnas.org/doi/abs/10.1073/pnas.1813263115> doi:
1016 10.1073/pnas.1813263115
- 1017 Myhre, G., Alterskjær, K., Stjern, C. W., Hodnebrog, Ø., Marelle, L., Samset, B. H.,
1018 ... Stohl, A. (2019, Nov 05). Frequency of extreme precipitation increases
1019 extensively with event rareness under global warming. *Scientific Reports*, 9(1),
1020 16063.
- 1021 Naveau, P., Hannart, A., & Ribes, A. (2020). Statistical Methods for Extreme Event
1022 Attribution in Climate Science. *Annual Review of Statistics and Its Appli-*
1023 *cation*, 7(1), 89-110. Retrieved from [https://doi.org/10.1146/annurev-
1024 -statistics-031219-041314](https://doi.org/10.1146/annurev-statistics-031219-041314) doi: 10.1146/annurev-statistics-031219-041314
- 1025 Norwood, A., Kalnay, E., Ide, K., Yang, S.-C., & Wolfe, C. (2013, Jun). Lyapunov,
1026 singular and bred vectors in a multi-scale system: an empirical exploration
1027 of vectors related to instabilities. *Journal of Physics A: Mathematical and*
1028 *Theoretical*, 46(25), 254021. Retrieved from [https://dx.doi.org/10.1088/
1029 1751-8113/46/25/254021](https://dx.doi.org/10.1088/1751-8113/46/25/254021) doi: 10.1088/1751-8113/46/25/254021
- 1030 O'Brien, T. A., Collins, W. D., Kashinath, K., Rübel, O., Byna, S., Gu, J., ... Ull-
1031 rich, P. A. (2016). Resolution dependence of precipitation statistical fidelity in
1032 hindcast simulations. *Journal of Advances in Modeling Earth Systems*, 8(2),
1033 976-990. Retrieved from [https://agupubs.onlinelibrary.wiley.com/doi/
1034 abs/10.1002/2016MS000671](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2016MS000671) doi: <https://doi.org/10.1002/2016MS000671>
- 1035 O'Gorman, P. A. (2015, Jun 01). Precipitation Extremes Under Climate Change.
1036 *Current Climate Change Reports*, 1(2), 49-59. Retrieved from <https://>

- doi.org/10.1007/s40641-015-0009-3 doi: 10.1007/s40641-015-0009-3
- 1037 Palmer, T. N., Buizza, R., Doblas-Reyes, F., Jung, T., Leutbecher, M., Shutts,
1038 G. J., ... Weisheimer, A. (2009). Stochastic parametrization and model
1039 uncertainty. *ECMWF Technical Memoranda*.
- 1040 Palmer, T. N., & Zanna, L. (2013, Jun). Singular vectors, predictability and en-
1041 semble forecasting for weather and climate. *Journal of Physics A: Mathemati-
1042 cal and Theoretical*, 46(25), 254018. Retrieved from [https://dx.doi.org/10](https://dx.doi.org/10.1088/1751-8113/46/25/254018)
1043 .1088/1751-8113/46/25/254018 doi: 10.1088/1751-8113/46/25/254018
- 1044 Pavliotis, G. A. (2014). *Stochastic processes and applications: diffusion processes,
1045 the Fokker-Planck and Langevin equations* (Vol. 60). Springer.
- 1046 Pazo, D., Rodriguez, M. A., & Lopez, J. M. (2010). Spatio-temporal evolution of
1047 perturbations in ensembles initialized by bred, Lyapunov and singular vectors.
1048 *Tellus A*, 62(1), 10-23. Retrieved from [https://onlinelibrary.wiley.com/
1049 doi/abs/10.1111/j.1600-0870.2009.00419.x](https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1600-0870.2009.00419.x) doi: [https://doi.org/10.1111/
1050 j.1600-0870.2009.00419.x](https://doi.org/10.1111/j.1600-0870.2009.00419.x)
- 1051 Pfahl, S., O’Gorman, P. A., & Fischer, E. M. (2017, Jun 01). Understanding the
1052 regional pattern of projected future changes in extreme precipitation. *Nature
1053 Climate Change*, 7(6), 423-427. Retrieved from [https://doi.org/10.1038/
1054 nclimate3287](https://doi.org/10.1038/nclimate3287) doi: 10.1038/nclimate3287
- 1055 Qi, D., & Majda, A. J. (2016). Predicting fat-tailed intermittent probability distri-
1056 butions in passive scalar turbulence with imperfect models through empirical
1057 information theory. *Communications in Mathematical Sciences*, 14(6), 1687-
1058 1722.
- 1059 Ragone, F., & Bouchet, F. (2021). Rare Event Algorithm Study of Extreme Warm
1060 Summers and Heatwaves Over Europe. *Geophysical Research Letters*, 48(12),
1061 e2020GL091197. Retrieved from [https://agupubs.onlinelibrary.wiley
1062 .com/doi/abs/10.1029/2020GL091197](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2020GL091197) (e2020GL091197 2020GL091197) doi:
1063 <https://doi.org/10.1029/2020GL091197>
- 1064 Ragone, F., Wouters, J., & Bouchet, F. (2018). Computation of extreme
1065 heat waves in climate models using a large deviation algorithm. *Proceed-
1066 ings of the National Academy of Sciences*, 115(1), 24-29. Retrieved from
1067 <https://www.pnas.org/content/115/1/24> doi: 10.1073/pnas.1712645115
- 1068 Saha, A., & Ravela, S. (2022). *Downscaling Extreme Rainfall Using Physical-
1069 Statistical Generative Adversarial Learning*. Retrieved from [https://
1070 arxiv.org/abs/2212.01446](https://arxiv.org/abs/2212.01446)
- 1071 Sapsis, T. P. (2020). Output-weighted optimal sampling for Bayesian regression
1072 and rare event statistics using few samples. *Proceedings of the Royal Society
1073 A: Mathematical, Physical and Engineering Sciences*, 476(2234), 20190834.
1074 Retrieved from [https://royalsocietypublishing.org/doi/abs/10.1098/
1075 rspa.2019.0834](https://royalsocietypublishing.org/doi/abs/10.1098/rspa.2019.0834) doi: 10.1098/rspa.2019.0834
- 1076 Schmidli, J., Goodess, C. M., Frei, C., Haylock, M. R., HunDecha, Y., Ribalaygua,
1077 J., & Schmith, T. (2007). Statistical and dynamical downscaling of precipita-
1078 tion: An evaluation and comparison of scenarios for the European Alps. *Jour-
1079 nal of Geophysical Research: Atmospheres*, 112(D4). Retrieved from [https://
1080 agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2005JD007026](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2005JD007026) doi:
1081 <https://doi.org/10.1029/2005JD007026>
- 1082 Schorlepp, T., Tong, S., Grafke, T., & Stadler, G. (2023, October). Scal-
1083 able methods for computing sharp extreme event probabilities in infinite-
1084 dimensional stochastic systems. *Statistics and Computing*, 33(6). Re-
1085 trieved from <http://dx.doi.org/10.1007/s11222-023-10307-2> doi:
1086 10.1007/s11222-023-10307-2
- 1087 Sheshadri, A., Borrus, M., Yoder, M., & Robinson, T. (2021). Midlatitude Error
1088 Growth in Atmospheric GCMs: The Role of Eddy Growth Rate. *Geophys-
1089 ical Research Letters*, 48(23), e2021GL096126. Retrieved from [https://
1090 agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2021GL096126](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2021GL096126)
- 1091

- 1092 (e2021GL096126 2021GL096126) doi: <https://doi.org/10.1029/2021GL096126>
1093 Sterk, A. E., & van Kekem, D. L. (2017, Sep 24). Predictability of Extreme Waves
1094 in the Lorenz-96 Model Near Intermittency and Quasi-Periodicity. *Complex-*
1095 *ity, 2017*, 9419024. Retrieved from <https://doi.org/10.1155/2017/9419024>
1096 doi: 10.1155/2017/9419024
- 1097 Tandon, N. F., Zhang, X., & Sobel, A. H. (2018). Understanding the Dy-
1098 namics of Future Changes in Extreme Precipitation Intensity. *Geo-*
1099 *physical Research Letters, 45*(6), 2870-2878. Retrieved from [https://](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2017GL076361)
1100 agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2017GL076361 doi:
1101 <https://doi.org/10.1002/2017GL076361>
- 1102 Tantet, A., van der Burgt, F. R., & Dijkstra, H. A. (2015). An early warning in-
1103 dicator for atmospheric blocking events using transfer operators. *Chaos: An*
1104 *Interdisciplinary Journal of Nonlinear Science, 25*(3), 036406. Retrieved from
1105 <https://doi.org/10.1063/1.4908174> doi: 10.1063/1.4908174
- 1106 Thompson, V., Dunstone, N. J., Scaife, A. A., Smith, D. M., Slingo, J. M.,
1107 Brown, S., & Belcher, S. E. (2017, Jul 24). High risk of unprecedented
1108 UK rainfall in the current climate. *Nature Communications, 8*(1), 107.
1109 Retrieved from <https://doi.org/10.1038/s41467-017-00275-3> doi:
1110 10.1038/s41467-017-00275-3
- 1111 Touchette, H. (2009). The large deviation approach to statistical mechanics. *Physics*
1112 *Reports, 478*(1-3), 1–69.
- 1113 Uribe, F., Papaioannou, I., Marzouk, Y. M., & Straub, D. (2021). Cross-Entropy-
1114 Based Importance Sampling with Failure-Informed Dimension Reduction for
1115 Rare Event Simulation. *SIAM/ASA Journal on Uncertainty Quantification,*
1116 *9*(2), 818-847. Retrieved from <https://doi.org/10.1137/20M1344585> doi:
1117 10.1137/20M1344585
- 1118 van der Wiel, K., Kapnick, S. B., Vecchi, G. A., Cooke, W. F., Delworth, T. L.,
1119 Jia, L., ... Zeng, F. (2016). The Resolution Dependence of Contiguous U.S.
1120 Precipitation Extremes in Response to CO2 Forcing. *Journal of Climate,*
1121 *29*(22), 7991 - 8012. Retrieved from [https://journals.ametsoc.org/](https://journals.ametsoc.org/view/journals/clim/29/22/jcli-d-16-0307.1.xml)
1122 [view/journals/clim/29/22/jcli-d-16-0307.1.xml](https://journals/clim/29/22/jcli-d-16-0307.1.xml) doi: 10.1175/
1123 JCLI-D-16-0307.1
- 1124 Villén-Altamirano, M., Villén-Altamirano, J., et al. (1991). RESTART: a method
1125 for accelerating rare event simulations. *Queueing, Performance and Control in*
1126 *ATM (ITC-13)*, 71–76.
- 1127 Walter R. Gilks, G. O. R., & Sahu, S. K. (1998). Adaptive Markov Chain Monte
1128 Carlo through Regeneration. *Journal of the American Statistical Associ-*
1129 *ation, 93*(443), 1045-1054. Retrieved from [https://doi.org/10.1080/](https://doi.org/10.1080/01621459.1998.10473766)
1130 [01621459.1998.10473766](https://doi.org/10.1080/01621459.1998.10473766) doi: 10.1080/01621459.1998.10473766
- 1131 Wang, Q., Mu, M., & Sun, G. (2020, Jan 01). A useful approach to sensitivity and
1132 predictability studies in geophysical fluid dynamics: conditional non-linear
1133 optimal perturbation. *National Science Review, 7*(1), 214-223. Retrieved from
1134 <https://doi.org/10.1093/nsr/nwz039> doi: 10.1093/nsr/nwz039
- 1135 Wasserman, L. (2004). *All of statistics*. New York: Springer New York, NY. doi: 10
1136 .1007/978-0-387-21736-9
- 1137 Webber, R. J., Plotkin, D. A., O'Neill, M. E., Abbot, D. S., & Weare, J. (2019).
1138 Practical rare event sampling for extreme mesoscale weather. *Chaos: An In-*
1139 *terdisciplinary Journal of Nonlinear Science, 29*(5), 053109. Retrieved from
1140 <https://doi.org/10.1063/1.5081461> doi: 10.1063/1.5081461
- 1141 Wilks, D. S. (2005). Effects of stochastic parametrizations in the Lorenz '96 system.
1142 *Quarterly Journal of the Royal Meteorological Society, 131*(606), 389-407. Re-
1143 trieved from [https://rmets.onlinelibrary.wiley.com/doi/abs/10.1256/](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1256/qj.04.03)
1144 [qj.04.03](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1256/qj.04.03) doi: <https://doi.org/10.1256/qj.04.03>
- 1145 Wouters, J., & Bouchet, F. (2016, Aug). Rare event computation in deter-
1146 ministic chaotic systems using genealogical particle analysis. *Journal*

- 1147 *of Physics A: Mathematical and Theoretical*, 49(37), 374002. Retrieved
 1148 from <https://dx.doi.org/10.1088/1751-8113/49/37/374002> doi:
 1149 10.1088/1751-8113/49/37/374002
- 1150 Wouters, J., Schiemann, R. K. H., & Shaffrey, L. C. (2023). Rare Event Simu-
 1151 lation of Extreme European Winter Rainfall in an Intermediate Complexity
 1152 Climate Model. *Journal of Advances in Modeling Earth Systems*, 15(4),
 1153 e2022MS003537. Retrieved from [https://agupubs.onlinelibrary.wiley](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2022MS003537)
 1154 [.com/doi/abs/10.1029/2022MS003537](https://doi.org/10.1029/2022MS003537) (e2022MS003537 2022MS003537) doi:
 1155 <https://doi.org/10.1029/2022MS003537>
- 1156 Wright, D. B., SaMaras, C., & Lopez-Cantu, T. (2021). Resilience to Extreme
 1157 Rainfall Starts with Science. *Bulletin of the American Meteorological Society*,
 1158 102(4), E808 - E813. Retrieved from [https://journals.ametsoc.org/
 1159 view/journals/bams/102/4/BAMS-D-20-0267.1.xml](https://journals.ametsoc.org/view/journals/bams/102/4/BAMS-D-20-0267.1.xml) doi: 10.1175/
 1160 BAMS-D-20-0267.1
- 1161 Yiou, P., & Jezequel, A. (2020). Simulation of extreme heat waves with empirical
 1162 importance sampling. *Geoscientific Model Development*, 13(2), 763–781. Re-
 1163 trieved from <https://gmd.copernicus.org/articles/13/763/2020/> doi: 10
 1164 .5194/gmd-13-763-2020
- 1165 Zhang, B. J., Sahai, T., & Marzouk, Y. M. (2022). A Koopman framework for rare
 1166 event simulation in stochastic differential equations. *Journal of Computational
 1167 Physics*, 456, 111025. Retrieved from [https://www.sciencedirect.com/
 1168 science/article/pii/S0021999122000870](https://www.sciencedirect.com/science/article/pii/S0021999122000870) doi: [https://doi.org/10.1016/
 1169 j.jcp.2022.111025](https://doi.org/10.1016/j.jcp.2022.111025)
- 1170 Zuckerman, D. M., & Chong, L. T. (2017). Weighted Ensemble Simula-
 1171 tion: Review of Methodology, Applications, and Software. *Annual Re-
 1172 view of Biophysics*, 46(1), 43-57. Retrieved from [https://doi.org/
 1173 10.1146/annurev-biophys-070816-033834](https://doi.org/10.1146/annurev-biophys-070816-033834) (PMID: 28301772) doi:
 1174 10.1146/annurev-biophys-070816-033834
- 1175 Zuev, K. (2015). *Subset Simulation Method for Rare Event Estimation: An Intro-
 1176 duction*. Retrieved from <https://arxiv.org/abs/1505.03506>