

1 Pseudo-prospective Evaluation of
2 UCERF3-ETAS Forecasts During the 2019
3 Ridgecrest Sequence

4

5 William H. Savran, Maximilian J. Werner, Warner Marzocchi, David A. Rhoades, David D.

6 Jackson, Kevin Milner, Edward Field, Andrew Michael

7

8 William Savran, University of Southern California, Southern California Earthquake Center, 3651

9 Trousdale Parkway, Los Angeles, CA 90089

10

11 Abstract

12 The 2019 Ridgecrest sequence provides the first opportunity to evaluate Uniform California
13 Earthquake Rupture Forecast Version 3 with Epidemic Type Aftershock Sequences (UCERF3-
14 ETAS) in a pseudo-prospective sense. For comparison, we include a version of the model
15 without explicit faults more closely mimicking traditional ETAS models (UCERF3-NoFaults).
16 We evaluate the forecasts with new metrics developed within the Collaboratory for the Study of
17 Earthquake Predictability (CSEP). The metrics consider synthetic catalogs simulated by the
18 models rather than synoptic probability maps, thereby relaxing the Poisson assumption of
19 previous CSEP tests. Our approach compares statistics from the synthetic catalogs directly
20 against observations, providing a flexible approach that can account for dependencies and
21 uncertainties encoded in the models. We find that, to first order, both UCERF3-ETAS and
22 UCERF3-NoFaults approximately capture the spatiotemporal evolution of the Ridgecrest
23 sequence, adding to the growing body of evidence that ETAS models can be informative
24 forecasting tools. However, we also find that both models mildly overpredict the seismicity rate,
25 on average, aggregated over the evaluation period. More severe testing indicates the
26 overpredictions occur too often for observations to be statistically indistinguishable from the
27 model. Magnitude tests indicate that the models do not include enough variability in forecasted
28 magnitude-number distributions to match the data. Spatial tests highlight discrepancies between
29 the forecasts and observations, but the greatest differences between the two models appear when
30 aftershocks occur on modeled UCERF3-ETAS faults. Therefore, any predictability associated
31 with embedding earthquake triggering on the (modeled) fault network may only crystalize during
32 the presumably rare sequences with aftershocks on these faults. Accounting for uncertainty in the
33 model parameters could improve test results during future experiments.

34 Introduction

35 A fundamental question in seismology is: What is the probability of observing an earthquake
36 within some predefined space-time-magnitude region? Earthquake forecasting models try to
37 answer this question by incorporating ideas of varying complexity about the earthquake process,
38 including both empirical statistical relations, such as the Omori-Utsu and Gutenberg-Richter
39 relations (Gutenberg and Richter, 1944; Utsu, 1961), and physical modeling, such as Coulomb
40 stress calculations (Oppenheimer et al., 1988; King et al., 1994; Stein, 1999; Woessner et al.,
41 2011; Cattania et al., 2018). The simplest models use locations of previous earthquakes to
42 forecast locations of future earthquakes via smoothing (Kagan and Jackson, 1994; Frankel, 1995;
43 Werner et al., 2010; Zechar and Jordan, 2010; Werner et al., 2011; Helmstetter and Werner,
44 2014). By contrast, UCERF3-ETAS (hereafter U3ETAS) combines long-term earthquake
45 probabilities on faults based on elastic rebound statistics with short-term earthquake clustering as
46 epidemic type aftershock sequences (Ogata, 1998) into a single model with fault-specific
47 magnitude distributions (Field et al., 2017a; Field et al., 2017b). Most notably, U3ETAS
48 provides probabilities of triggering ruptures on known faults, such as the Garlock and San
49 Andreas faults. U3ETAS is a candidate model for Operational Earthquake Forecasting (OEF)
50 issued by the US Geological Survey, motivating model evaluations also from a practical
51 perspective.

52 The Collaboratory for the Study of Earthquake Predictability (CSEP) has established the
53 philosophy and cyber-infrastructure required to conduct earthquake forecasting experiments in
54 an unbiased and transparent fashion (Jordan, 2006; Schorlemmer and Gerstenberger, 2007;
55 Jordan et al., 2011; Michael and Werner, 2018; Schorlemmer et al., 2018). Since its inception,
56 CSEP has been using likelihood-based consistency tests (Schorlemmer et al., 2007; Zechar et al.,

57 2010; Rhoades et al., 2011; Werner et al., 2011) that are rooted in the concepts that 1)
58 earthquakes occur in space-time-magnitude bins independently, 2) earthquakes follow the
59 Poisson distribution in each bin, and 3) modelers provide the 'true' parameter of the distribution
60 in each bin. Thus, CSEP required that modelers provide forecasts giving the expected number of
61 earthquakes in discrete space-time-magnitude bins. This pragmatic simplification allows multiple
62 types of models, including those without explicit likelihood functions, to participate in the
63 experiments.

64 However, Poisson likelihood-based evaluations of gridded forecasts can incorrectly
65 report discrepancies between forecasts and observations when the true likelihood function of a
66 forecast does not match a Poisson distribution or when strong dependencies exist between events
67 within a forecast period. For example, the ETAS model is overdispersed with respect to a
68 Poisson process, causing forecasts to be more frequently rejected than expected (Werner and
69 Sornette 2008, Lombardi and Marzocchi 2010, Nandan et al., 2019). This is particularly
70 noticeable when evaluating forecasts over multiple forecasting periods.

71 Evaluating gridded forecasts over multiple time periods exploits the property that the sum
72 of N Poisson random variables each with parameter λ_i is a Poisson random variable with
73 parameter $\sum \lambda_i$. The same convenience does not hold for catalog-based forecasts, because, in
74 general, simulated events in catalogs from later time periods are not consistent with simulated
75 events from earlier catalogs. Thus, catalog-based forecasting models should be evaluated for
76 consistency by comparing realizations from their predictive distributions against observations.
77 This approach is formally referred to as calibration, which is based on the idea that observations
78 should be indistinguishable from realizations drawn from the predictive distributions of the
79 model (Gneiting et al., 2006; Gneiting et al., 2007; Gneiting and Katzfuss, 2014). In other words,

80 if the model were the data generator, we would expect observations to uniformly sample the
81 forecasted distribution over independent trials.

82 Fundamentally, calibration is a different type of evaluation approach that can be
83 potentially more severe than previously used cumulative evaluations. For example, evaluations
84 over individual periods might indicate that observations consistently fall within the forecasted
85 distribution, but instead of sampling the forecasted distribution uniformly they are concentrated
86 towards one end. Thus, the model would fail calibration, but potentially pass a cumulative test.
87 Understanding the overall performance of these models is more important than ‘rejecting’ a
88 particular forecast; therefore, we focus on characteristics of the models and differences between
89 models that potentially uncover new insights that might lead to model improvements.

90 Page and van der Elst (2018) introduced Turing-style evaluations for assessing
91 forecasting models that produce synthetic catalogs. The tests evaluate important features of the
92 simulated catalogs such as: aftershock productivity, seismicity rate, magnitude distribution, and
93 clustering behavior. The Turing tests provide useful insights into the behavior of the forecasts,
94 and can help to inform modeling decisions and identify discrepancies between the model and
95 observations. However, they are not well suited for consistency testing or calibration, because
96 they do not formally score forecasts against observations.

97 Here, we introduce new validation methods (consistency tests) for catalog-based
98 earthquake forecasting models. Most notably, these methods relax the assumption that
99 earthquakes follow independent Poisson distributions in discrete space-time-magnitude bins
100 (Schorlemmer et al., 2007). Catalog-based forecasts differ from gridded forecasts in that they can
101 capture the full aleatory variability of the model and can also account for epistemic uncertainty
102 (such as in parameter estimates). Exhaustive sets of simulated catalogs retain the full

103 spatiotemporal dependencies amongst modeled earthquakes, i.e., they can reflect the full
104 complexity of the model through simulations. We build predictive distributions from the
105 forecasts, empirically, by defining statistics that emphasize important characteristics of
106 seismicity. This enables hypothesis testing and calibration of probabilistic forecasts over multiple
107 evaluation periods.

108 We organize this manuscript as follows. First, we introduce the evaluation metrics for
109 catalog-based forecasts. We then apply the metrics to forecasts made during the Ridgecrest
110 sequence for an eleven-week period following the M_w 7.1 mainshock. To benchmark the fault-
111 based triggering component of U3ETAS, we also generate and evaluate forecasts from a simpler
112 version of the model, named UCERF3-NoFaults (hereafter NoFaults), which removes the fault
113 component of U3ETAS. We discuss the primary differences between U3ETAS and NoFaults in
114 the Methods section. Finally, we discuss the evaluation results with respect to U3ETAS and
115 NoFaults and comment on the evaluation metrics.

116 [Methods: Evaluations](#)

117 [Definitions and Notation](#)

118 We introduce some notation to help us define evaluations in the context of earthquake
119 forecasts that are specified as synthetic earthquake catalogs. First, we define a testing region \mathcal{R} ,
120 as the combination of a magnitude range \mathcal{M} , spatial domain \mathcal{S} , and time period \mathcal{T} :

121

$$\mathcal{R} = \mathcal{M} \times \mathcal{S} \times \mathcal{T}. \quad (1)$$

122

123 These individual components can be regarded as filters that operate on a catalog which retain
124 only the events within \mathcal{R} .

125 Let us consider an event, $e = (t, \mathbf{x}, m)$. Each e can be specified exactly by its origin time,
126 t , geographic location, \mathbf{x} , and magnitude, m . The spatial coordinate, \mathbf{x} , typically refers to a
127 latitude and longitude pair, but can also include depth. Thus, an earthquake catalog is simply a
128 collection of events.

129 We define an observed catalog as

130

$$\Omega = \{e_i \mid i = 1, \dots, N_{obs}; e_i \in \mathcal{R}\}. \quad (2)$$

131

132 Here, Ω is the observed catalog containing N_{obs} observed events, e_i , within \mathcal{R} . This catalog is
133 used as the testing data set for the evaluations.

134 A forecast is a collection of synthetic catalogs whose events \tilde{e}_{ij} in \mathcal{R} are defined as

135

$$\Lambda \equiv \Lambda_j = \{\tilde{e}_{ij} \mid i = 1, \dots, N_j; j = 1, \dots, J; \tilde{e}_{ij} \in \mathcal{R}\}. \quad (3)$$

136

137 The forecast, Λ , contains J synthetic catalogs each with N_j events. Λ_j indicates the j^{th}
138 catalog of the forecast Λ , likewise \tilde{e}_{ij} denotes the i^{th} event from the j^{th} synthetic catalog of
139 Λ . Each Λ_j is a synthetic catalog that represents a continuous space-time-magnitude realization of
140 seismicity generated by the model. The synthetic catalogs from the forecast and the observed
141 catalog share the same event definitions, therefore the same statistics can be readily applied to all
142 catalogs.

143 The testing methodology presented here follows three guiding principles: (1) statistics
144 should be calculated directly from the simulated and observed catalogs to build test distributions
145 empirically; (2) testing methods should be able to preserve space-time-magnitude dependencies
146 between events that are encoded in the model and may exist within the earthquake process; and
147 (3) these tests should reduce their reliance on approximate likelihood functions of models,
148 whether parametric in the case of the Poisson assumption or non-parametric in the case of the
149 spatial test and pseudo-likelihood tests presented here. The last principle requires compromise if
150 (approximate) likelihood-based inference remains desirable for model comparison, especially if
151 no analytical likelihood function is available. Models without explicit likelihood functions are
152 also known as generative or simulator-based models (Gutmann and Corander, 2016), which is
153 the case for U3ETAS. In the remainder of this section, we define a suite of evaluations that can
154 be used to evaluate the consistency of earthquake forecasts specified as synthetic catalogs against
155 observed seismicity. These evaluations by no means represent an exhaustive set of metrics that
156 can be used to evaluate catalog-based forecast models.

157

158 [Number Test](#)

159 The number test asks whether the number of earthquakes observed in \mathcal{R} is inconsistent with the
160 forecasted number distribution by assessing whether the observed number falls into the tails of
161 the forecast distribution (Kagan and Jackson, 1995; Schorlemmer et al., 2007; Zechar et al.,
162 2010). The test statistic for an arbitrary catalog, ξ , is $N = |\xi|$, where the bars denote the count of
163 events in the catalog. Thus, the observed statistic is

164

$$N_{obs} = |\Omega|, \tag{4}$$

165 or simply the number of events in the observed catalog. To build the test distribution from the
166 forecast Λ we calculate this statistic for every catalog forming the vector:

167

$$N_j = |\Lambda_j|; j = 1, \dots, J. \quad (5)$$

168

169 To identify potentially important discrepancies between the observation and the forecast
170 distribution, we compute the quantiles of the observed number in the empirical cumulative
171 distribution function (CDF) of the forecast distribution (Equation 5) according to

172

$$\delta_1 = 1 - F_N(N_{obs} - 1) = P(N_j \geq N_{obs}) \quad (6)$$

173

174 and

$$\gamma_N = \delta_2 = F_N(N_{obs}) = P(N_j \leq N_{obs}). \quad (7)$$

175

176 $F_N(n)$ denotes the empirical cumulative distribution function of N_j . For the number test, we
177 should consider a two-sided test to assess the probabilities of observing (1) at least and (2) at
178 most N_{obs} events, a distinction that becomes important when forecasted and observed numbers
179 are small (Zechar et al., 2010). $F_N(n)$ denotes the empirical predictive CDF of N_j . For a
180 probabilistically calibrated forecast, we expect the quantile scores, γ_N , to uniformly sample the
181 forecasted number distribution over multiple independent trials.

182

183 Magnitude Test

184 The magnitude test evaluates whether an observed magnitude-frequency distribution (MFD) is
185 inconsistent with the forecasted MFD. We base this statistic on a square metric computed from
186 the difference in logarithms between the incremental MFDs of the so-called union catalog Λ_U ,
187 individual catalogs Λ_j , and the observed catalog Ω . This metric is loosely related to the quadratic
188 Cramer von-Mises and Anderson tests (Anderson, 2006). Using the logarithm of bin-wise
189 magnitude counts places greater weight on magnitude bins with relatively fewer observed (and
190 predicted) earthquakes, which typically occur at larger magnitudes. Thus, each missed (or over-
191 predicted) event at larger magnitudes should contribute more to the test statistic than the same
192 absolute error between smaller magnitudes.

193 We first define the union catalog Λ_U as

194

$$\Lambda_U = \{\Lambda_1 \cup \Lambda_2 \cup \dots \cup \Lambda_J\}. \quad (8)$$

195

196 The union catalog Λ_U contains all events from $\mathbf{\Lambda}$ totaling $N_U = \sum_{j=1}^J |\Lambda_j|$ events. We compute
197 the standard histograms of (1) $\Lambda_U^{(m)}$, the magnitudes of the union catalog, (2) $\Lambda_j^{(m)}$, the
198 magnitudes of each individual synthetic catalog, and (3) $\Omega^{(m)}$, the observed magnitudes, with all
199 histograms discretized according to \mathcal{M} (say, in increments of 0.1 magnitude units). We
200 normalize all histograms so that $\sum_k \xi^{(m)}(k) = N_{obs}$, where $\xi^{(m)}(k)$ represents the normalized
201 number of events in the k^{th} bin of the incremental MFD for an arbitrary catalog. This ensures
202 that differences in forecasted rates do not contribute directly to the bin-wise sum, although the
203 earthquake rate may implicitly affect the shape of the MFD. We compute the observed statistic

204 as the sum of squared logarithmic residuals between the normalized observed magnitude and
 205 union histograms following

206

$$d_{obs} = \sum_k \left(\log \left[\frac{N_{obs}}{N_U} \Lambda_U^{(m)}(k) + 1 \right] - \log [\Omega^{(m)}(k) + 1] \right)^2. \quad (9)$$

207

208 $\Lambda_U^{(m)}(k)$ and $\Omega^{(m)}(k)$ represent the count in the k^{th} bin of the incremental MFDs from the union
 209 and observed catalogs, respectively. We add unity to each bin to prevent the singularity
 210 associated with $\log(0)$. Since we are only concerned with differences between two MFDs, this
 211 modification does not bias the statistic. Next, we build the test distribution from the catalogs in
 212 Λ , i.e., the distribution of test statistics if the forecast model were the data-generating model
 213 following

214

$$D_j = \sum_k \left(\log \left[\frac{N_{obs}}{N_U} \Lambda_U^{(m)}(k) + 1 \right] - \log \left[\frac{N_{obs}}{N_j} \Lambda_j^{(m)}(k) + 1 \right] \right)^2 ; j = 1, \dots, J. \quad (10)$$

215

216 Here, $\Lambda_j^{(m)}(k)$ indicates the count of events in the k^{th} magnitude bin from the j^{th} synthetic
 217 catalog. Finally, we compute the quantile score of d_{obs} within the empirical cumulative
 218 distribution function defined as

219

$$\gamma_m = F_D(d_{obs}) = P(D_j \leq d_{obs}). \quad (11)$$

220

221 We expect the quantile scores, γ_m , should uniformly sample the test distribution D_j for either
222 forecast.

223

224 Pseudo-Likelihood Test

225 Here, we introduce a statistic based on the continuous point-process likelihood function (Daley
226 and Vere-Jones, 2004). While this statistic resembles the likelihood scores used by previous
227 CSEP experiments (e.g., Schorlemmer et al., 2007), there are two differences. First, we do not
228 compute an actual likelihood, whence the name pseudo-likelihood. Second, this pseudo-
229 likelihood statistic is aggregated over target event likelihood scores as opposed to the Poisson
230 likelihood scores computed over discrete cells (see also Rhoades et al., 2011). In the case of zero
231 or one events the pseudo-likelihood and the Poisson likelihood scores are identical. Finally, and
232 most importantly, we build test distributions of pseudo-likelihood scores using the simulated
233 (non-Poissonian) catalogs provided by the forecasting model, thereby producing distributions
234 that better represent models that are over-dispersed and more clustered than a Poisson process.

235 A continuous marked space-time point process can be represented by its conditional
236 intensify function $\lambda(\mathbf{e} | H_t)$, where H_t denotes the history of all earthquake occurrences (and any
237 other relevant input data) prior to time t . The log likelihood function of any point-process over a
238 region \mathcal{R} is

239

$$L = \sum_{i=1}^N \ln \lambda(e_i | H_t) - \int_{\mathcal{R}} \lambda(\mathbf{e} | H_t) d\mathcal{R}. \quad (12)$$

240

241 CSEP seeks to accommodate a wide range of stochastic models, including generative or
 242 simulator-based models such as UCERF3-ETAS without explicit conditional intensity or
 243 likelihood functions. CSEP therefore does not require an explicit likelihood function for
 244 evaluation (although models that contain explicit likelihood functions can be evaluated using this
 245 idea, e.g., Ogata et al., 2013).

246 Instead, we approximate the expectation of $\lambda(\mathbf{e} | H_t)$ using the forecasted catalogs. To do
 247 this we introduce a discretization of \mathcal{R} similar to previous CSEP experiments. Heuristically, the
 248 approximate rate density is defined as the conditional expectation, given the discretized region,
 249 \mathcal{R}_d , of its continuous rate density:

$$\hat{\lambda}(\mathbf{e} | H_t) = E[\lambda(\mathbf{e} | H_t) | \mathcal{R}_d]. \quad (13)$$

251
 252 Conceptually, we can still regard the model as continuous in space, time and magnitude,
 253 but its rate density is only approximated and takes a constant value within a given cell. The
 254 approximate rate density is readily derived from the standard CSEP forecast of gridded expected
 255 rates, by computing the mean event count from the forecast, $\mathbf{\Lambda}$, in each cell in \mathcal{R}_d . The discrete
 256 grid cells are used only for approximation purposes; we use the synthetic catalogs from the full
 257 model to calculate the pseudo-likelihood statistic (rather than catalogs of the approximate
 258 model).

259 From the approximate rate density (Equation 13), we can define the pseudo log likelihood
 260 \hat{L} by

$$\hat{L} = \sum_{i=1}^N \ln \hat{\lambda}(e_i | H_t) - \int_{\mathcal{R}} \hat{\lambda}(\mathbf{e} | H_t) d\mathcal{R}. \quad (14)$$

261 The pseudo-likelihood test applied here considers a discretized region in space to avoid
 262 introducing artifacts into the forecasts (such as minimum “water-levels” and smoothing operators
 263 that could bias the evaluations) to account for under-sampling in space-magnitude bins.
 264 Formally, we can write the spatial approximate rate density as

$$\hat{\lambda}_s(\mathbf{e} | H_t) = \sum_{\mathcal{M}} \hat{\lambda}(\mathbf{e} | H_t). \quad (15)$$

266
 267 If $\hat{\lambda}_s(k)$ denotes the approximate rate density in the k^{th} spatial cell of the model, we can
 268 compute the observed pseudo-likelihood score using,

$$\hat{L}_{obs} = \sum_{i=1}^{N_{obs}} \ln \hat{\lambda}_s(k_i) - \bar{N}. \quad (16)$$

270
 271 Here k_i denotes the spatial cell in which the i^{th} event occurs and \bar{N} denotes the expected number
 272 of events in \mathcal{R}_d . Following Equation 16, we compute the statistics for the test distribution as

$$\hat{L}_j = \left[\sum_{i=1}^{N_j} \ln \hat{\lambda}_s(k_{ij}) - \bar{N} \right]; j = 1, \dots, J. \quad (17)$$

274
 275 Here $\hat{\lambda}_s(k_{ij})$ denotes the approximate rate density of the i^{th} event of the j^{th} catalog from the
 276 forecast. We combine Equation 16 and Equation 17 to obtain the quantile score

277

$$\gamma_L = F_L(\hat{L}_{obs}) = P(\hat{L}_j \leq \hat{L}_{obs}). \quad (18)$$

278

279 The statistic captures simultaneously the spatial component and the rate component of the
 280 forecast. Thus, potential discrepancies in both rate and the spatial components of the forecasts
 281 should be reflected in this statistic. As with the magnitude test and the number test, we expect
 282 that the quantile scores γ_L should be uniformly distributed over multiple evaluation periods.

283

284 Spatial Test – Geometric Average of Target Event Rate Distribution

285 The spatial test isolates the spatial distribution of the forecast to evaluate whether the observed
 286 locations are consistent with the forecasted spatial distribution. This statistic utilizes the
 287 approximate rate density (Equation 15) with normalization $\hat{\lambda}_s^* = \hat{\lambda}_s / \sum_{\mathcal{R}} \hat{\lambda}_s$ to isolate the spatial
 288 component of the forecast.

289 We define the observed spatial statistic according to

290

$$S_{obs} = \left[\sum_{i=1}^{N_{obs}} \ln \hat{\lambda}_s^*(k_i) \right] N_{obs}^{-1}, \quad (19)$$

291

292 where $\hat{\lambda}_s^*(k_i)$ denotes the normalized approximate rate density in the k^{th} cell corresponding to
 293 the i^{th} event in Ω . Likewise, we can define the test distribution for the statistic defined in
 294 Equation (19) using

295

$$S_j = \left[\sum_{i=1}^{N_j} \ln \hat{\lambda}_s^*(k_{ij}) \right] N_j^{-1}; j = 1, \dots, J. \quad (20)$$

296

297 As above, $\hat{\lambda}_s^*(k_{ij})$ denotes the approximate rate density in the k^{th} cell corresponding to the i^{th}

298 event in the j^{th} simulated catalog. The observed spatial statistic (Equation 19) is scored by

299 computing quantiles in the test distribution (Equation 20) using,

300

$$\gamma_S = F_S(\hat{s}_{obs}) = P(\hat{S}_j \leq \hat{s}_{obs}). \quad (21)$$

301

302 We interpret this statistic as being the geometric mean of the target event rate

303 distribution. Normalizing $\hat{\lambda}_s$ and computing the geometric mean of the target event rate

304 distribution ensures that two catalogs (from the same forecast) with events occurring in identical

305 bins will result in equivalent spatial test statistics irrespective of the number of events in either

306 catalog. If the model were the data generator, we expect that γ_S will be uniformly distributed

307 over multiple evaluation periods.

308

309 Testing Over Multiple Periods

310 To assess models over many periods, we exploit the following idea: quantile scores over

311 multiple periods should be uniformly distributed if the model is the data generator (Gneiting and

312 Katzfuss, 2014). Departures from a uniform distribution of the quantile scores flag discrepancies

313 between the forecasting model and observation. Formally, we employ a Kolmogorov-Smirnov

314 test between the quantile scores and the uniform distribution to test the hypothesis that the

315 observed quantile scores are uniformly distributed. We calculate the p -value of this test and use a
316 significance level $\alpha = 0.05$ to identify discrepancies.

317 Graphically, we consider different patterns of variation of the observed quantile scores
318 from a uniform distribution. A model that under-predicts the test statistic produces a graph
319 similar to that in Figure 1a. In this case, there is a small proportion of low quantile scores and a
320 high proportion of high quantiles compared to the uniform distribution, because the observed
321 test-statistic tends to be higher than the simulated test-statistic. Conversely, a model that tends to
322 over-predict the test statistic produces a graph similar to Figure 1b, because in that case the
323 actual test statistic tends to be lower than the simulated test statistics. If the model test statistics
324 are under-dispersed relative to the observed test statistics, then the quantile scores will fall near
325 the end-points 0 and 1 of the distribution. This produces the pattern seen in Figure 1c.
326 Conversely, if the model test statistics are over-dispersed relative to the actual test statistic, the
327 pattern seen in Figure 1d will be the result.

328

329 [Methods: Pseudo-Prospective Experiment Design](#)

330 The 2019 Ridgecrest sequence provides the first opportunity to evaluate operational aftershock
331 forecasts in a pseudo-prospective sense. A pseudo-prospective experiment preserves the time-
332 dependent causality of the data set by partitioning the dataset into a training set and a testing set
333 (Schorlemmer et al., 2018), which happened naturally as these forecasts were computed in near
334 real-time during the Ridgecrest sequence. Most of the forecasts produced in this study were
335 computed in near-real-time using real-time data products with the exceptions listed in Table 1.
336 The forecasts presented in this study use the *ShakeMap (v.14)* source model and default

337 parameters described in Milner et al. (2020). We evaluate the forecasts starting at $t = 0$ and $t =$
338 7 days following the M_w 7.1 mainshock (along with the nine others) in this study.

339

340 Data

341 For this experiment, we use authoritative data from the Advanced National Seismic
342 System (ANSS) provided by the United States Geological Survey (USGS) Comprehensive
343 Catalog (ComCat). The evaluation data were accessed from ComCat on 11 November 2019,
344 approximately 60 days following the date of the final forecast. We use the data directly provided
345 by ComCat, and do not attempt to standardize magnitude types or manually relocate events. We
346 apply the time-dependent magnitude of completeness model from Helmstetter et al. (2006) to
347 account for missing events following the mainshock, modeled by a time-dependent magnitude of
348 completeness $M_c(t)$. Therefore, the evaluation catalog has a threshold magnitude

349

$$M_t(t) = \max(M_{min}, M_c(t)). \quad (22)$$

350

351 Here, M_{min} , represents a minimum magnitude that is either defined to be $M_{min} = 2.5$ or $M_{min} =$
352 3.5, in the case of the number test. We apply the time-dependent magnitude of completeness
353 model to both forecasted and observed catalogs. The inset in Figure 2a shows the events used for
354 this study along with the time-dependent magnitude of completeness. In the 77 days following
355 the M_w 7.1 Ridgecrest event, the catalog lists 1,362 events with $M \geq 2.5$ in the study region.

356 Finite-fault representations for trigger ruptures are based on surface field mapping and
357 geodetic observations, and were provided by ShakeMap (Wald et al., 1999). These finite-fault
358 models were made available on 11 July 2019 within six days after the M_w 7.1 mainshock. Milner

359 et al. (2020) explains the various finite-fault representations available and the sensitivity of the
360 forecasts to these.

361

362 Earthquake Forecasting Models

363 We consider two forecasting models, (1) UCERF3-ETAS (U3ETAS) and (2) UCERF3-
364 NoFaults (NoFaults). The former model is explained in detail by Field et al. (2017b), so we
365 summarize the important differences between U3ETAS and NoFaults here. Field et al. (2017a)
366 provides a less technical overview of the UCERF3-ETAS model for the interested reader. The
367 full mathematical description of these models can be found in the above manuscripts and their
368 appendices.

369 U3ETAS is unique as compared with standard ETAS models, because the model includes
370 finite faults that can host so-called suprasedismogenic earthquakes. In U3ETAS, a
371 suprasedismogenic earthquake is defined as an earthquake with a rupture length at least as long as
372 the seismogenic fault width. When a large earthquake in close enough proximity to a U3ETAS
373 fault is sampled by ETAS, that earthquake is mapped onto the modeled fault-sections.
374 Subsequently, the rates of all events that utilize the ruptured sections are modified according to
375 Reid renewal statistics (Reid, 1910; Field et al., 2015). Therefore, U3ETAS provides stochastic
376 event sets with ruptures on modeled finite faults in addition to ‘off-fault’ ruptures elsewhere,
377 following a traditional ETAS model. U3ETAS makes no model-wide assumptions about
378 magnitude-frequency distributions on faults, with most exhibiting non Gutenberg-Richter (GR)
379 behavior depending on the relative rate of microseismicity versus inferred fault-based ruptures.
380 On average the faults are slightly characteristic (elevated rates at higher magnitudes), which
381 means off-fault areas are slightly anti-characteristic so that combined a GR b -value of 1.0 is

382 maintained. However, the model assigns the regional faults surrounding the Ridgecrest sequence
383 an anti-characteristic behavior (Field et al., 2017b; Milner et al., 2020) implying lower
384 probabilities of triggering supraseismogenic aftershocks than under a pure GR model. In
385 contrast, NoFaults applies the state-wide G-R relationship (b -value=1.0) throughout the entire
386 model.

387 As the name suggests, NoFaults does not include information about modeled faults and
388 behaves similar to a traditional space-time ETAS implementation (Ogata and Zhuang, 2006).
389 Both U3ETAS and NoFaults explicitly model the depth distribution of seismicity. The
390 computational requirements for the two models differ by approximately an order of magnitude,
391 with U3ETAS being more expensive. Because model simplicity and computational efficiency are
392 two desirable characteristics of robust operational forecasting tools (Jordan and Jones, 2010;
393 Jordan et al., 2011), we seek to understand the relative predictive skills and usefulness of the
394 models.

395 The forecasts issued by both U3ETAS and NoFaults consist of a family of 100,000
396 synthetic catalogs constrained to the bounding-box of the CSEP California testing region
397 (Schorlemmer and Gerstenberger, 2007). As inputs to all forecasts, we include earthquakes with
398 $M_{\text{w}} 2.5+$ for seven days prior to the $M_{\text{w}} 7.1$ mainshock until the start-time of each forecast,
399 including the $M_{\text{w}} 6.4$ Searles Valley event. We use identical input catalogs for both U3ETAS and
400 NoFaults to maintain direct comparability between the two forecasts. Also, we do not include
401 spontaneous (background) events in the conditioning data for the forecasts. Therefore, any
402 discrepancies between the forecast and observations can be attributed to the implementation of
403 the short-term components of the model and not the background seismicity model.

404

405 Spatial Region, Magnitude Bins, and Forecast Horizons

406 For this experiment, we choose magnitude bins

407

$$\mathcal{M} = \{[2.5, 2.6), [2.7, 2.8), \dots, [8.4, 8.5), [8.5, \infty)\}. \quad (23)$$

408

409 The bins are uniformly spaced at $\Delta M = 0.1$ except for the right-most bin which extends
410 to infinity. We remove events outside a spatial zone of three Wells and Coppersmith (1994) fault
411 radii from the M7.1 epicenter (143 km) to isolate the Ridgecrest aftershocks from other
412 seismicity. Each forecast horizon extends for seven non-overlapping days, which we treat as
413 independent time intervals. Figure 2b shows the spatial extent of the circular region surrounding
414 the hypocenter of the M_w 7.1 mainshock and the observed M2.5+ events during this time period.

415 These definitions completely define the extent of \mathcal{R} for our experiment. All forecasts are
416 evaluated for seven days following the forecast start time to preserve effects of short-term
417 clustering in the observed catalog. Table 1 contains the exact start and end times for all the
418 forecasts considered in this study, which consist of eleven non-overlapping time periods
419 following the M_w 7.1 mainshock. All but two U3ETAS forecasts were computed prospectively
420 using real-time catalogs and data. The NoFaults simulations were run pseudo-prospectively, but
421 using the same input catalogs and input finite-fault models as U3ETAS.

422 Results

423 Before we share the results of the quantitative evaluations of the forecasts, we show how
424 differences between U3ETAS and NoFaults manifest in individual synthetic catalogs. Since the
425 models are similar for events smaller than $\sim M_w$ 6.5, catalogs display similar characteristics for

426 ‘typical’ realizations (Figure 3a,c), defined here as catalogs representing the median of the
427 forecasted number distribution. The differences become obvious when viewing catalogs (Figure
428 3b,d) that sample the tails of the number distribution at the 99.9th percentile. We call these
429 catalogs ‘extreme’ as they forecast rare, but possible, large aftershock sequences on potentially
430 multiple faults. Extreme U3ETAS scenarios involve ruptures triggered on the Garlock fault and
431 subsequently on the San Andreas fault. Their respective aftershocks are largely constrained
432 within the fault zones. On the other hand, NoFaults assigns aftershock locations isotropically in
433 space resulting in (nearly) isotropic catalogs that contain clusters of earthquakes (Figure 3d).

434 The differences illustrated in Figure 3, namely in the catalogs at the tails of the forecast,
435 complicate robust model comparisons using typical California aftershock sequences, which only
436 occasionally involve triggering of large aftershocks on (mapped) faults. This is because the
437 models produce very similar (visually nearly indistinguishable) catalogs near the modes and
438 medians of the number distributions. Sequences such as the 1992 Landers earthquake cascade
439 and others that are thought to have triggered other large ruptures could help distinguish between
440 these two models (Kisslinger and Jones, 1991; Hauksson et al., 1993; Freed and Lin, 2001).

441 We show test results as quantile scores for all evaluations in Table 2. The overall
442 (aggregate) scores over all forecast periods are reported as p -values of Kolmogorov-Smirnov
443 tests between a uniform distribution and the quantile scores of each test computed at the updating
444 periods shown in Table 1.

445

446 Forecasted Seismicity Rates

447 Figure 4 shows the forecasted number distributions as a function of time during the
448 aftershock sequence for both $M_t(t) = \max(2.5, M_c(t))$ and $M_t(t) = \max(3.5, M_c(t))$.

449 We observe the largest variability in the number distribution immediately following the
450 mainshock, which decreases rapidly throughout the evaluation period. During the first evaluation
451 period the median forecasted numbers are 925 and 956 for U3ETAS and NoFaults, respectively,
452 with 829 observed events during this period. The median forecasted event counts are identical
453 between the two models for the remaining forecasting periods.

454 We compute number test results for each forecast by reporting quantile scores for
455 individual testing periods as a function of evaluation day (Figure 5a). Except for the first day,
456 both forecasts produce nearly identical quantile scores. The difference in number distributions
457 during the first forecasting period can potentially be explained by the anti-characteristic behavior
458 of the U3ETAS faults surrounding the aftershock sequence. This behavior results in U3ETAS
459 producing fewer large (M6.5+) events, along with their numerous aftershocks, and subsequently
460 fewer catalogs with large numbers of aftershocks. During the first forecasting period, the 95-
461 percentile range of the number distribution are (751, 2482) and (756, 3906) for U3ETAS and
462 NoFaults respectively.

463 Figure 5b shows the number test quantile scores compared against standard uniform
464 quantiles as a quantile-quantile plot. We assign the standard uniform quantiles following $U^{(k)} =$
465 $k/(n + 1)$, for $k = 1, \dots, n$, to space the quantiles equally along the distribution. We compute
466 confidence intervals for the k^{th} order statistic of the standard uniform distribution using $U^{(k)} \sim$
467 $B(k, n + 1 - k)$ where n is the number of observations (Jones, 2004).

468 The distribution of quantile scores indicates the forecasts overpredict the observed
469 seismicity (Figure 1b), as the observed numbers of earthquakes too frequently fall into the lower
470 tails of the forecasted distributions. At both magnitude cutoffs, the Kolmogorov-Smirnov tests
471 reject the hypothesis that the distribution of quantile scores from the number test are uniformly

472 distributed. This suggests that, given this limited forecasting period, the observations are not
473 indistinguishable from realizations from the forecast number distribution.

474

475 Magnitude Number Distribution

476 Figure 6a shows incremental MFDs aggregated over the eleven-week evaluation period.

477 For the union MFD, $\Lambda_U^{(m)}$, and observed MFDs, $\Omega_U^{(m)}$, we sum bin-wise counts from each

478 evaluation period to obtain aggregate counts. We estimate percentiles using an aggregate

479 forecasted MFD (thin lines in Figure 6). We generate the aggregate forecasted MFD using a

480 bootstrapped approach where we randomly sample one MFD per forecast per time-period and

481 sum bin-wise counts over each evaluation period. This produces 100,000 aggregate MFDs

482 approximating an MFD representative of the eleven-week evaluation period. Except between

483 M3.0 and M4.0 the observations generally fall within the variability of the forecasted MFD.

484 Above M6.5 we see differences in the tails of the magnitude frequency distributions that further

485 show how the anticharacteristic MFDs assumed by U3ETAS manifest in the forecasts.

486 Figure 6b shows the bin-wise value of the magnitude test statistic over the full evaluation

487 period to highlight the bin-wise contribution to the overall magnitude test statistic. From the bin-

488 wise statistic, we can obtain the magnitude test statistic in Equation 9 by summing over all

489 magnitude bins. This figure illustrates that discrepancies at larger magnitudes contribute more

490 (per event) to the value magnitude test statistic, but this must be reconciled with statistics

491 computed from simulated catalogs. We can identify bins whose values contribute most to the

492 discrepancy between observations and forecasts by assessing the observed statistic with respect

493 to the bin-wise distribution of magnitude test statistics.

494 The percentiles in Figure 6b (for both U3ETAS and NoFaults) are estimated from the
495 bin-wise distributions of magnitude test statistics based on the bootstrapped aggregate MFD
496 (explained above). We can associate the large peak observed near M4.7 in Figure 6b with
497 catalogs from either model that contain zero events in that magnitude bin. This can be seen by
498 comparing the square bin-wise difference with the union MFD and zero observed events in
499 Figure 6a. The percentiles in Figure 6b indicate 2.5% of the catalogs contain no events at this
500 magnitude, and 16% of the catalogs contain no events at M5.0. The largest discrepancies with
501 respect to the forecasts occur from around M3.0 through M4.0 indicated by the observed bin-
502 wise values falling outside the 95th percentile range of the bin-wise distribution. Generally, the
503 observed values are frequently greater than the median from their respective bin-wise
504 distributions, and this behavior is not confined to a particular magnitude range.

505 Figure 7a shows quantile scores for each evaluation period following the M_w 7.1
506 mainshock to assess the performance of the forecast over multiple updating periods. The shaded
507 region in Figure 7a indicates the critical region assuming a 0.05 significance level for a right-
508 tailed statistical test. (In this magnitude test, larger-than-expected values of the statistic, i.e. large
509 quantile scores, indicate larger discrepancies). Figure 7b shows the quantile-quantile plot of the
510 magnitude test scores against standard uniform quantiles. The quantile scores, γ_m , do not sample
511 the test distribution uniformly and are instead concentrated near the upper end. The Kolmogorov-
512 Smirnov test thus rejects the hypothesis of a uniform distribution of the quantile scores. The
513 pattern in Figure 7b implies persistently greater-than-expected differences between the observed
514 magnitude distribution and the forecast. The pattern of magnitude test quantile scores reflects the
515 finding in Figure 6b that the bin-wise magnitude scores are typically greater than the median bin-
516 wise values.

517 Spatial Distribution of Seismicity and Pseudo-likelihood Test

518 Figure 8a,b shows the approximate spatial rate density (Equation 15) for both U3ETAS
519 and NoFaults during the first evaluation period following the M_w 7.1 mainshock. The expected
520 cell-wise event counts clearly show differences between U3ETAS and NoFaults, specifically the
521 increased expected rates along modeled faults in U3ETAS. The relatively high rates along the
522 Garlock fault, for example, are dominated by catalogs containing suprasedismogenic ruptures
523 along these faults (which occur in about 7% of the catalogs). Thus, we should expect to see
524 noticeable differences between these two models with observations of such aftershock
525 sequences.

526 Figure 8c shows test distributions of spatial statistics for a single week-long forecast
527 immediately following the M_w 7.1 mainshock. Likewise, Figure 8d shows test distributions for
528 the pseudo-likelihood score. Positive values of the pseudo-likelihood scores can occur when
529 multiple target events occur within the same spatial bin with $\hat{\lambda}_s \gg 1$ (the Poisson likelihood
530 contains an explicit term to account this discretization artifact that does not appear in the pseudo-
531 likelihood statistic), which can happen when scoring catalogs that sample upper tails of the
532 number distribution. For this evaluation period, the observed statistic, \hat{L}_{obs} , lies in the lower tail
533 of the test distribution \hat{L} .

534 The aggregate spatial test result in Figure 9a shows quantile scores and pseudo-
535 likelihood quantiles for each evaluation period since the M_w 7.1 mainshock. In general, U3ETAS
536 tends to have larger quantile scores, and thus, more favorable test statistics for a given forecast
537 than NoFaults. We find that if differences are observed, they appear in both the pseudo-
538 likelihood and spatial test statistics. Comparisons of quantile scores against the uniform
539 distribution (Figure 9b) show the statistic from the observed catalog tends to fall in the lower tail

540 of the spatial test distribution for most forecasts. Thus, according to the spatial test, random
541 draws from the forecasted distribution are distinguishable from the observations; the latter more
542 frequently occur in cells of lower rates than expected by the models.

543 The pseudo-likelihood quantiles γ_L shows seemingly better agreement with the standard
544 uniform quantiles (we compute $p=0.0280$, $p=0.0235$ from the Kolmogorov-Smirnov test for
545 U3ETAS and NoFaults, respectively); however, this observation must be analyzed in the context
546 of both the number test and spatial test results. Since both models show inconsistencies in the
547 number test and spatial test, we expect this to be reflected in the pseudo-likelihood test. Previous
548 studies have shown that the Poisson-based likelihood test is anticorrelated with the number test
549 results (Werner et al., 2011). The somewhat counterintuitive result causes forecasts that
550 overpredict the seismicity rates to trivially pass the likelihood test. Therefore, this must be
551 considered when interpreting the pseudo-likelihood test results. Specifically, the test results are
552 probably better solely because the models overpredict.

553 Deconstructing the statistics helps to inform us about the behavior of the evaluation
554 results. For the magnitude test, we showed the bin-wise value of the test statistics to identify
555 problematic bins. Here, we show cell-wise spatial pseudo-likelihood ratios (U3ETAS –
556 NoFaults) in Figure 10 to understand which cells contribute to the differences observed in the
557 spatial test and the pseudo-likelihood tests. We represent the observed event rate distribution on
558 the spatial grid as follows: spatial cells with no observed events show the difference in the
559 approximate rate density between models, and cells containing observed events show the
560 difference in the that cells' contribution to the pseudo-likelihood scores. Only cells containing
561 observed events contribute to the spatial test statistic, therefore cells without observed events
562 help to visualize differences in the spatial distributions of the forecast. These plots are similar to

563 spatial deviance residuals (Schneider et al., 2014). We find that U3ETAS tends to show larger
564 spatial test statistics, and thus quantile scores, when observed events occur along modeled
565 U3ETAS faults.

566 Discussion

567 We have introduced a suite of evaluations for catalog-based earthquake forecasts that provide
568 insight into the forecasted earthquake rates, magnitude-number distributions, and spatial
569 distributions of seismicity. These evaluations are complementary to the Turing Tests introduced
570 by Page and van der Elst (2018) and the comparative mean-information gain introduced by
571 Nandan et al. (2019), which can also be used to evaluate generative or simulator models that
572 produce synthetic catalogs. Importantly, these metrics begin to relax the independence and
573 Poisson assumptions of previous forecast evaluations (Schorlemmer et al., 2007). Additionally,
574 we introduced an approach, commonly applied to weather (and other) probabilistic forecasts
575 (e.g., Gneiting et al., 2006), to calibrate probabilistic earthquake forecasting models. We apply
576 these new methods to U3ETAS and NoFaults forecasts of the Ridgecrest sequence for eleven-
577 weeks following the Mw 7.1 mainshock.

578 We find U3ETAS and NoFaults overpredict earthquake rates in 10 out of 11 evaluation
579 periods for $M_t(t) = \max(2.5, M_c(t))$ by comparing observed event counts against the mode of
580 the forecasted number distribution (modal ratio), but 5 out of 11 modal ratios are within $\pm 20\%$ of
581 the observed event count (with the maximum being a 140% overprediction). On average, from
582 the modal ratio, the forecasts overpredict observed event counts by approximately 50%.
583 NoFaults tends to produce larger variability in the number distribution than U3ETAS (e.g.,
584 Figure 4a,b), which is most noticeable during the first evaluation period. This likely occurs
585 because the Airport Lake and Little Lake faults are both anti-characteristic in U3ETAS (Milner

586 et al., 2020), which causes these faults to host fewer large magnitude events as compared with
587 the GR ($b=1.0$) MFD implemented in NoFaults (Figure 6). Moreover, every event in NoFaults is
588 treated as a point-source. In contrast, U3ETAS can assign large ruptures to faults (if the event
589 occurs close enough to a modeled fault). This in turn activates the elastic-rebound model (Field
590 et al., 2015), and this combined behavior effectively smooths the forecasted number of events in
591 the vicinity of the rupture (Figure 8a,b). The anticharacteristic behavior of the Little Lake and
592 Airport Lake faults is likely to have pronounced differences in the tails of the number
593 distributions and the associated hazard and risk curves. In areas with anticharacteristic MFDs,
594 U3ETAS produces lower expected rates of events except along the faults that host aftershock
595 sequences. Visually, we see the larger rates along the faults for U3ETAS as compared with
596 NoFaults (Figure 8a,b), but statistically the chance of damaging aftershocks is lower in U3ETAS.

597 On aggregate, the U3ETAS and NoFaults produce catalogs whose MFDs display lower
598 variability with respect to the expected MFD than observations. By comparing the logarithms of
599 bin-wise counts we find that observations are different, statistically, from realizations from the
600 forecasts. Figure 6b shows contributions to this discrepancy across all magnitude ranges, but
601 M3.0 through M4.0 show the largest discrepancy with respect to the forecasted bin-wise
602 statistics. This can be interpreted in two ways: either significant discrepancies exist between
603 U3ETAS (and NoFaults) and observations, or this magnitude test is too severe given the
604 uncertainties in reported magnitudes and assumed b -values in the forecasting model. To address
605 uncertainties in reported magnitudes, we recomputed the magnitude test with magnitude bins
606 $\Delta M = 0.2$, and found consistent results with those presented in Figure 7. Moreover, using Monte
607 Carlo simulations we find the magnitude test results are sensitive to changes in b -values of $\Delta b \leq$
608 0.1 units, which is on the order of the uncertainty in b -value estimates for U3ETAS (Felzer,

609 2013). Thus, including epistemic uncertainty in the assumed b -value could potentially improve
610 calibration. Furthermore, we should consider explicitly accounting for uncertainties in observed
611 magnitudes when evaluating earthquake forecasts.

612 Here, we discuss a potential reason for the inconsistencies in the spatial test results.
613 ETAS models, due to their self-excitation property (Hawkes, 1971), have a particularly difficult
614 time forecasting seismicity in areas that were not previously active. As a result, the approximate
615 rate densities (Equation 13) and locations of events in the simulated catalogs are controlled by
616 the events in the input catalog used to condition the forecast. For example, neither U3ETAS nor
617 NoFaults forecast much seismicity off the northwest-end of the mainshock fault plane during the
618 first forecasting period (Figure 8a,b), leading to the observations falling in the lower tail of the
619 test distribution. This discrepancy can be reduced with more frequent updating of the ETAS
620 intensity function, which would locally increase after each subsequent event. Ideally, the
621 conditional intensity function would be updated continuously after each observed event;
622 however, this might prove difficult in practice because of computational times and costs.

623 The spatial and pseudo-likelihood tests show the largest differences between U3ETAS
624 and NoFaults amongst the statistics, which we expected because the spatial distribution of
625 seismicity is the primary difference between these models. Figure 10 shows spatial (pseudo-)
626 log-likelihood ratios (U3ETAS – NoFaults) to understand where differences in the spatial test
627 statistic originate. Carefully looking at the cell-wise ratios where observed events occur in Figure
628 10, we find the differences manifest when aftershocks occur near modeled U3ETAS faults. This
629 suggests that we should be able to identify differences between U3ETAS and NoFaults using the
630 spatial test for sequences when aftershocks occur on modeled U3ETAS faults.

631 We draw counter-intuitive conclusions from the pseudo-likelihood test, when put in
632 context of the spatial and number tests. We find that observations are inconsistent with both the
633 rate and spatial forecasts from both models, and thus we expect the pseudo-likelihood scores to
634 reflect this observation. Instead, the pseudo-likelihood test scores show more favorable
635 agreement with the observations. Similar to the Poisson likelihood test (Schorlemmer et al.,
636 2007), overpredictions in rates can result in artificially high pseudo-likelihood scores (e.g.,
637 Werner et al., 2011). From this, we conclude that the pseudo-likelihood test provides redundant
638 information to the number and spatial tests, and the test is less severe than the spatial test when
639 the forecast fails the number test.

640 U3ETAS uses ETAS parameters estimated from the state-wide California seismic catalog
641 (Hardebeck, 2013). The moderate overprediction by U3ETAS (and NoFaults) suggests that the
642 Ridgecrest sequence deviates from the state-wide average in aftershock productivity. Milner et
643 al. (2020) found this behavior was due to high primary productivity of the mainshock, coupled
644 with low secondary aftershock productivity. State-wide maximum-likelihood estimates (MLE) of
645 ETAS parameters also result in over-predictions for this sequence when using traditional ETAS
646 models (Mancini et al., 2020, In Press). These results suggest that accurate forecasting of
647 aftershock rates requires proper treatment of intersequence variability or obtaining sequence
648 specific parameters (Page et al., 2016).

649 MLE parameter estimates of a traditional ETAS model may well be different, however,
650 from MLE estimates of U3ETAS parameters, because the models are different: non-GR behavior
651 in U3ETAS is spatially variable, magnitude and spatial distributions are not separable, and
652 ‘characteristic-ness’ impacts secondary triggering productivity (Milner et al., 2020). Milner et al.
653 (2020) showed that adjustment of the ETAS c -value could improve the fit to the cumulative

654 number of $M \geq 3.5$ events, but this required manual trial-and-error adjustments to optimize for a
655 specific metric. If sequence specific parameters are not (yet) available, incorporating additional
656 uncertainty in the ETAS parameters could make the model more general and perhaps calibrated,
657 especially for the first forecasts following a large earthquake before sequence-specific
658 information is available (Omi et al., 2015; Omi et al., 2019).

659 The discrepancies between the models and observations can potentially be explained by
660 epistemic uncertainty in model parameters not accounted for by the model. Incorporating
661 parameter uncertainty would broaden distribution functions (reduce sharpness) and potentially
662 lead to calibrated probabilistic forecasts. Moreover, incorporating more sequences (and quiet
663 periods) could uncover systematic discrepancies with observations that can lead to improvements
664 in the models, and increase the robustness of the results. Retrospective as well as further
665 prospective tests are required to understand the usefulness and accuracy of modeling decisions.
666 In particular, the U3ETAS model will be most easily differentiated from standard ETAS models
667 in the rare circumstances (of about 7%, assuming U3ETAS is correct) when suprasedismogenic
668 events are triggered. This relatively small percentage (which varies spatially in the model)
669 implies that we expect to observe substantial differences between the models about once in 20
670 earthquake sequences. Future work should therefore evaluate the model retrospectively against
671 all well-recorded aftershock sequences observed in California.

672 Conclusions

673 In this manuscript, we evaluate forecasts from UCERF3-ETAS and UCERF3-NoFaults during
674 the Ridgecrest using new non-parametric evaluations developed for forecasts specified as
675 simulated catalogs. We evaluate eleven week-long forecasts immediately following the Mw 7.1
676 mainshock using an idea, known as calibration, that suggests that random draws from the

677 forecast should be indistinguishable from observations. Probabilistic calibration is severe, but is
678 a useful approach to aggregate forecasts over multiple periods. Probabilistic forecasts should aim
679 to maximize the sharpness of their predictive distributions, subject to calibration (Gneiting et al.,
680 2007; Gneiting and Katzfuss, 2014). We introduce statistics that probe the forecasted earthquake
681 rate, magnitude distributions, and spatial component of the forecast. Importantly, these
682 evaluations relax the assumption that earthquakes occur in discrete Poissonian space-time-
683 magnitude bins and better reflect the dependencies between earthquakes.

684 This pseudo-prospective evaluation of U3ETAS (and NoFaults) constitutes a milestone as
685 it represents the first out-of-sample evaluation of a model under consideration for real-time
686 operational earthquake forecasting by the US Geological Survey. (Pseudo-) Prospective model
687 evaluation is a critical step in building confidence in the model outputs. To first order, both
688 U3ETAS and NoFaults capture the temporal evolution and magnitude distributions of the
689 earthquake sequence, notwithstanding the generic state-wide ETAS model parameters. For
690 example, when considering the mode of the forecasted number distribution, the forecasts on
691 average overpredict the observed number of events by approximately 50% with 5 out of 11
692 forecasts being within $\pm 20\%$ of the observed event count. This suggests that, in spite of the much
693 more severe calibration test results, U3ETAS (and ETAS models in general) are effective tools
694 to provide insight into the spatial and temporal distributions of seismicity, in real-time, during an
695 aftershock sequence. As with any forecasting model, the usefulness depends on the specific use-
696 case in mind (Field and Milner, 2018). For U3ETAS, in particular, estimates of probabilities of
697 ruptures on nearby faults may provide valuable information for emergency planners and decision
698 makers (Milner et al., 2020).

699 The results of the proposed tests lead to similar conclusions for both UCERF3-ETAS and
700 NoFaults. For the number test, the forecasts systematically overpredict the observed seismicity.
701 The overpredictions can be attributed to deviations in primary and secondary aftershock
702 productivity during the Ridgecrest with respect to the state wide average. The observed MFDs
703 show greater variability with respect to the expected MFD than predicted by the forecasts. We
704 interpret this discrepancy as a result of unmodeled uncertainty in the magnitude data,
705 highlighting a need to account for observational uncertainty in the tests. The spatial tests uncover
706 an issue associated with the discrete updating of self-exciting ETAS models, that is, the models
707 have difficulty forecasting seismicity in areas without previous seismicity. We find the largest
708 differences between U3ETAS and NoFaults when observed aftershocks occur on modeled
709 U3ETAS faults. In such cases, the pseudo-likelihood test provides redundant results to the
710 number and spatial test.

711 Data and Resources

712 The evaluation results and data for individual simulations can be found at
713 https://github.com/cseptesting/ridgecrest_evaluation_bssa. The UCERF3-ETAS and UCERF3-
714 NoFaults simulations were generated using the UCERF3 model implemented in OpenSHA and
715 can be found at <https://github.com/opensha/ucerf3-etlas-launcher/>. The code used for the analysis
716 can be found in development at <https://github.com/SCECcode/csep2/>. The finite-fault data was
717 obtained from the ShakeMap accessed through the Comprehensive Catalog (ComCat) provided
718 by the United States Geological Survey and can be access through the web at
719 <https://earthquake.usgs.gov/data/comcat/>.

720 Acknowledgements

721 We would like to thank Morgan Page, Jeanne Hardebeck, and Nicholas van der Elst for their
722 insight and helpful comments regarding forecast evaluations. M.J.W. and W.M. received funding
723 from the European Union's Horizon 2020 research and innovation program (No 821115, RISE:
724 Real-Time Earthquake Risk Reduction for a Resilient Europe). This research was supported by
725 the Southern California Earthquake Center (Contribution No. 10082). SCEC is funded by NSF
726 Cooperative Agreement EAR-1600087 & USGS Cooperative Agreement G17AC00047.

727 References

- 728 Anderson, T. W. (2006). On The Distribution of the Two-Sample Cramer von-Mises Criterion,
729 *Annals of Mathematical Statistics* 1-12.
- 730
- 731 Cattania, C., M. J. Werner, W. Marzocchi, S. Hainzl, D. Rhoades, M. Gerstenberger, M. Liukis,
732 W. Savran, A. Christophersen, A. Helmstetter, A. Jimenez, S. Steacy, and T. H. Jordan (2018).
733 The Forecasting Skill of Physics-Based Seismicity Models during the 2010-2012 Canterbury,
734 New Zealand, Earthquake Sequence, *Seismological Research Letters* **89** 1238-1250.
- 735
- 736 Daley, D. J., and D. Vere-Jones (2004). Scoring probability forecasts for point processes: the
737 entropy score and information gain, *Journal of Applied Probability* **41** 297-312.
- 738
- 739 Felzer, K. R. (2013). Appendix L: Estimate of the Seismicity Rate and Magnitude-Frequency
740 Distribution of Earthquakes in California from 1850 to 2011, 1-13.
- 741

742 Field, E. H., G. P. Biasi, P. Bird, T. E. Dawson, K. R. Felzer, D. D. Jackson, K. M. Johnson, T.
743 H. Jordan, C. Madden, A. J. Michael, K. R. Milner, M. T. Page, T. Parsons, P. M. Powers, B. E.
744 Shaw, W. R. Thatcher, R. J. Weldon II, and Y. Zeng (2015). Long-Term Time-Dependent
745 Probabilities for the Third Uniform California Earthquake Rupture Forecast (UCERF3), Bulletin
746 of the Seismological Society of America **105** 511-543.

747

748 Field, E. H., T. H. Jordan, M. T. Page, K. R. Milner, B. E. Shaw, T. E. Dawson, G. P. Biasi, T.
749 Parsons, J. L. Hardebeck, A. J. Michael, R. J. Weldon, P. M. Powers, K. M. Johnson, Y. H.
750 Zeng, K. R. Felzer, N. van der Elst, C. Madden, R. Arrowsmith, M. J. Werner, and W. R.
751 Thatcher (2017a). A Synoptic View of the Third Uniform California Earthquake Rupture
752 Forecast (UCERF3), Seismological Research Letters **88** 1259-1267.

753

754 Field, E. H., and K. R. Milner (2018). Candidate Products for Operational Earthquake
755 Forecasting Illustrated Using the HayWired Planning Scenario, Including One Very Quick (and
756 Not-So-Dirty) Hazard-Map Option, Seismological Research Letters **89** 1420-1434.

757

758 Field, E. H., K. R. Milner, J. L. Hardebeck, M. T. Page, N. J. van der Elst, T. H. Jordan, A. J.
759 Michael, B. E. Shaw, and M. J. Werner (2017b). A Spatiotemporal Clustering Model for the
760 Third Uniform California Earthquake Rupture Forecast (UCERF3-ETAS): Toward an
761 Operational Earthquake Forecast, Bulletin of the Seismological Society of America **107** 1049-
762 1081.

763

764 Frankel, A. (1995). Mapping Seismic Hazard in the Central and Eastern United States, **66** 8-21.

765

766 Freed, A. M., and J. Lin (2001). Delayed triggering of the 1999 Hector Mine earthquake by
767 viscoelastic stress transfer, *Nature* **411** 180-183.

768

769 Gneiting, T., F. Balabdaoui, and A. E. Raftery (2007). Probabilistic forecasts, calibration and
770 sharpness, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **69** 243-
771 268.

772

773 Gneiting, T., and M. Katzfuss (2014). Probabilistic forecasting, *Annual Review of Statistics and*
774 *Its Application* **1** 125-151.

775

776 Gneiting, T., K. Larson, K. Westrick, M. G. Genton, and E. Aldrich (2006). Calibrated
777 Probabilistic Forecasting at the Stateline Wind Energy Center, **101** 968-979.

778

779 Gutenberg, B., and C. F. Richter (1944). Frequency of earthquakes in California, *Bulletin of the*
780 *Seismological society of America* **34** 185-188.

781

782 Gutmann, M. U., and J. Corander (2016). Bayesian optimization for likelihood-free inference of
783 simulator-based statistical models, *The Journal of Machine Learning Research* **17** 4256-4302.

784

785 Hardebeck, J. L. (2013). Appendix S: Constraining Epidemic Type Aftershock Sequence
786 (ETAS) Parameters from the Uniform California Earthquake Rupture Forecast, Version 3

787 Catalog and Validating the ETAS Model for Magnitude 6.5 or Greater Earthquakes, U.S. Geol.
788 Surv. Open-File Rept.
789
790 Hauksson, E., L. M. Jones, K. Hutton, and D. Eberhart-Phillips (1993). The 1992 Landers
791 Earthquake Sequence: Seismological observations, **98** 19835.
792
793 Hawkes, A. G. (1971). Point spectra of some mutually exciting point processes, Journal of the
794 Royal Statistical Society: Series B (Methodological) **33** 438-443.
795
796 Helmstetter, A., Y. Y. Kagan, and D. D. Jackson (2006). Comparison of short-term and time-
797 independent earthquake forecast models for southern California, Bulletin of the Seismological
798 Society of America **96** 90-106.
799
800 Helmstetter, A., and M. J. Werner (2014). Adaptive Smoothing of Seismicity in Time, Space,
801 and Magnitude for Time-Dependent Earthquake Forecasts for California, Bulletin of the
802 Seismological Society of America **104** 809-822.
803
804 Jones, M. (2004). Families of distributions arising from distributions of order statistics, Test **13**
805 1-43.
806
807 Jordan, T. H. (2006). Earthquake predictability, brick by brick, pubs.geoscienceworld.org **77** 3-6.
808

809 Jordan, T. H., Y. T. Chen, P. Gasparini, R. Madariaga, I. Main, W. Marzocchi, and G.
810 Papadopoulos (2011). Operational earthquake forecasting. State of knowledge and guidelines for
811 utilization, *Annals of Geophysics*.
812

813 Jordan, T. H., and L. M. Jones (2010). Operational Earthquake Forecasting: Some Thoughts on
814 Why and How, *Seismological Research Letters* **81** 571-574.
815

816 Kagan, Y. Y., and D. D. Jackson (1994). Long-Term Probabilistic Forecasting of Earthquakes,
817 *Journal of Geophysical Research-Solid Earth* **99** 13685-13700.
818

819 King, G. C., R. S. Stein, and J. Lin (1994). Static stress changes and the triggering of
820 earthquakes, *Bulletin of the Seismological Society of America* **84** 935-953.
821

822 Kisslinger, C., and L. M. Jones (1991). Properties of aftershock sequences in southern California,
823 *Journal of Geophysical Research* **96** 11947.
824

825 Mancini, S., M. J. Werner, M. Segou, and T. Parsons (2020, In Press). The Predictive Skills of
826 Elastic Coulomb Rate-and-state Aftershock Forecasts During the 2019 Ridgecrest, California,
827 Earthquake Sequence, *Bulletin of the Seismological Society of America*.
828

829 Michael, A. J., and M. J. Werner (2018). Preface to the Focus Section on the Collaboratory for
830 the Study of Earthquake Predictability (CSEP): New Results and Future Directions,
831 *Seismological Research Letters* **89** 1226-1228.

832

833 Milner, K. R., E. H. Field, W. H. Savran, M. T. Page, and T. H. Jordan (2020). Operational
834 Earthquake Forecasting during the 2019 Ridgecrest, California, Earthquake Sequence with the
835 UCERF3-ETAS Model, *Seismological Research Letters*.

836

837 Nandan, S., G. Ouillon, D. Sornette, and S. Wiemer (2019). Forecasting the Full Distribution of
838 Earthquake Numbers Is Fair, Robust, and Better, *Seismological Research Letters*.

839

840 Ogata, Y. (1998). Space-time point-process models for earthquake occurrences, *Ann I Stat Math*
841 **50** 379-402.

842

843 Ogata, Y., K. Katsura, G. Falcone, K. Nanjo, and J. Zhuang (2013). Comprehensive and Topical
844 Evaluations of Earthquake Forecasts in Terms of Number, Time, Space, and Magnitude, *Bulletin*
845 *of the Seismological Society of America* **103** 1692-1708.

846

847 Ogata, Y., and J. Zhuang (2006). Space-time ETAS models and an improved extension,
848 *Tectonophysics* **413** 13-23.

849

850 Omi, T., Y. Ogata, Y. Hirata, and K. Aihara (2015). Intermediate-term forecasting of aftershocks
851 from an early aftershock sequence: Bayesian and ensemble forecasting approaches, *Journal of*
852 *Geophysical Research: Solid Earth* **120** 2561-2578.

853

854 Omi, T., Y. Ogata, K. Shiomi, B. Enescu, K. Sawazaki, and K. Aihara (2019). Implementation of
855 a Real-Time System for Automatic Aftershock Forecasting in Japan, *Seismological Research*
856 *Letters* **90** 242-250.

857

858 Oppenheimer, D. H., P. A. Reasenber, and R. W. Simpson (1988). Fault plane solutions for the
859 1984 Morgan Hill, California, earthquake sequence: Evidence for the state of stress on the
860 Calaveras fault, *Journal of Geophysical Research: Solid Earth* **93** 9007-9026.

861

862 Page, M. T., and N. J. van der Elst (2018). Turing-Style Tests for UCERF3 Synthetic Catalogs,
863 *Bulletin of the Seismological Society of America* **108** 729-741.

864

865 Page, M. T., N. J. van der Elst, J. Hardebeck, K. Felzer, and A. J. Michael (2016). Three
866 Ingredients for Improved Global Aftershock Forecasts: Tectonic Region, Time-Dependent
867 Catalog Incompleteness, and Intersequence Variability, *Bulletin of the Seismological Society of*
868 *America* **106** 2290-2301.

869

870 Reid, H. F. (1910). The California earthquake of April 18, 1906: Report of the State Earthquake
871 Investigation Commission. 2. The mechanics of the earthquake, Carnegie Inst. of Washington.

872

873 Rhoades, D. A., D. Schorlemmer, M. C. Gerstenberger, A. Christophersen, J. D. Zechar, and M.
874 Imoto (2011). Efficient testing of earthquake forecasting models, *Acta Geophys* **59** 728-747.

875

876 Schneider, M., R. Clements, D. A. Rhoades, and D. Schorlemmer (2014). Likelihood- and
877 residual-based evaluation of medium-term earthquake forecast models for California,
878 *Geophysical Journal International* **198** 1307-1318.

879

880 Schorlemmer, D., M. Gerstenberger, S. Wiemer, D. D. Jackson, and D. A. Rhoades (2007).
881 Earthquake likelihood model testing, *Seismological Research Letters* **78**.

882

883 Schorlemmer, D., and M. C. Gerstenberger (2007). RELM Testing Center, *Seismological*
884 *Research Letters* **78** 30-36.

885

886 Schorlemmer, D., M. J. Werner, W. Marzocchi, T. H. Jordan, Y. Ogata, D. D. Jackson, S. Mak,
887 D. A. Rhoades, M. C. Gerstenberger, N. Hirata, M. Liukis, P. J. Maechling, A. Strader, M.
888 Taroni, S. Wiemer, J. D. Zechar, and J. C. Zhuang (2018). The Collaboratory for the Study of
889 Earthquake Predictability: Achievements and Priorities, *Seismological Research Letters* **89** 1305-
890 1313.

891

892 Stein, R. S. (1999). The role of stress transfer in earthquake occurrence, *Nature* **402** 605-609.

893

894 Utsu, T. (1961). A statistical study on the occurrence of aftershocks, *Geophys. Mag.* **30** 521-605.

895

896 Wald, D. J., V. Quitoriano, T. H. Heaton, H. Kanamori, C. W. Scrivner, and C. B. Worden
897 (1999). TriNet “ShakeMaps”: Rapid generation of peak ground motion and intensity maps for
898 earthquakes in southern California, *Earthquake Spectra* **15** 537-555.

899

900 Wells, D. L., and K. J. Coppersmith (1994). New empirical relationships among magnitude,
901 rupture length, rupture width, rupture area, and surface displacement, *Bulletin of the*
902 *Seismological Society of America* **84** 974-1002.

903

904 Werner, M. J., A. Helmstetter, D. D. Jackson, and Y. Y. Kagan (2011). High-Resolution Long-
905 Term and Short-Term Earthquake Forecasts for California, *Bulletin of the Seismological Society*
906 *of America* **101** 1630-1648.

907

908 Werner, M. J., A. Helmstetter, D. D. Jackson, Y. Y. Kagan, and S. Wiemer (2010). Adaptively
909 smoothed seismicity earthquake forecasts for Italy, *Annals of Geophysics* **53** 107-116.

910

911 Woessner, J., S. Hainzl, W. Marzocchi, M. J. Werner, A. M. Lombardi, F. Catalli, B. Enescu, M.
912 Cocco, M. C. Gerstenberger, and S. Wiemer (2011). A retrospective comparative forecast test on
913 the 1992 Landers sequence, *Journal of Geophysical Research-Solid Earth* **116**.

914

915 Zechar, J. D., M. C. Gerstenberger, and D. A. Rhoades (2010). Likelihood-Based Tests for
916 Evaluating Space-Rate-Magnitude Earthquake Forecasts, *Bulletin of the Seismological Society*
917 *of America* **100** 1184-1195.

918

919 Zechar, J. D., and T. H. Jordan (2010). Simple smoothed seismicity earthquake forecasts for
920 Italy, *Annals of Geophysics* 1-7.

921

922 Tables

923 Table 1. Start times for the forecasts considered in this study. UCERF3-NoFaults were computed pseudo-
 924 prospectively using the same input catalogs as their UCERF3-ETAS counterparts. UCERF3-ETAS forecasts were
 925 computed in near-real-time with real-time data products except as otherwise noted.

Mw 7.1 + ΔT (days)	Start Time (GMT+0)	End Time (GMT+0)
0.0*	2019-07-06 03:19:54.04	2019-07-13 03:19:54.04
7.0†	2019-07-13 03:19:54.04	2019-07-20 03:19:54.04
14.0**	2019-07-20 03:19:54.04	2019-07-27 03:19:54.04
21.0	2019-07-27 03:19:54.04	2019-08-03 03:19:54.04
28.0	2019-08-03 03:19:54.04	2019-08-10 03:19:54.04
35.0	2019-08-10 03:19:54.04	2019-08-17 03:19:54.04
42.0	2019-08-17 03:19:54.04	2019-08-24 03:19:54.04
49.0	2019-08-24 03:19:54.04	2019-08-31 03:19:54.04
56.0	2019-08-31 03:19:54.04	2019-09-07 03:19:54.04
63.0	2019-09-07 03:19:54.04	2019-09-14 03:19:54.04
70.0	2019-09-14 03:19:54.04	2019-09-21 03:19:54.04

* Calculated on 09/04/19, catalog input data accessed from ComCat 09/04/19

† Calculated on 07/16/19, catalog input data accessed from ComCat 07/16/19

** Calculated on 08/19/19, catalog input data accessed from ComCat 08/19/19

926

927 Table 2. Evaluation results for number, magnitude, pseudo-likelihood, and spatial tests results for UCERF3-ETAS
 928 and UCERF3-NoFaults for $M_t(t) = \max(2.5, M_c(t))$.

Test day (since M7.1)	U3ETAS (N-Test)*	NoFaults (N-Test)*	U3ETAS (M-Test)	NoFaults (M-Test)	U3ETAS (PL-Test)	NoFaults (PL-Test)	U3ETAS (S-Test)	NoFaults (S-Test)
7	[0.818, 0.185]	[0.843, 0.160]	0.912	0.944	0.073	0.094	0.044	0.035
14	[0.688, 0.326]	[0.692, 0.322]	0.819	0.822	0.035	0.032	0.043	0.04
21	[0.995, 0.006]	[0.996, 0.006]	0.129	0.136	0.109	0.083	0.192	0.137
28	[0.958, 0.052]	[0.958, 0.051]	0.725	0.731	0.018	0.017	0.065	0.065
35	[0.999, 0.002]	[0.999, 0.002]	0.57	0.575	0.031	0.036	0.296	0.298
42	[0.907, 0.114]	[0.908, 0.113]	0.825	0.827	0.018	0.012	0.116	0.078
49	[0.399, 0.636]	[0.398, 0.636]	0.782	0.781	0.325	0.186	0.307	0.209
56	[0.998, 0.004]	[0.998, 0.004]	0.904	0.905	0.012	0.013	0.266	0.276
63	[0.999, 0.002]	[0.999, 0.002]	0.908	0.905	0.187	0.095	0.921	0.874
70	[0.995, 0.008]	[0.995, 0.008]	0.905	0.904	0.052	0.024	0.732	0.609
77	[1.000, 0.000]	[1.000, 0.001]	0.967	0.967	0.134	0.138	0.975	0.976
Overall	8.450e-05	3.363e-05	1.425e-03	1.222e-03	2.796e-02	2.349e-02	2.432e-06	1.927e-08

*(δ^1, δ^2)

929

930 Figure Captions

931 Figure 1. Schematic of cumulative distribution of quantile scores for a test statistic calculated
932 over multiple test periods (points) as compared with the ideal uniform distribution (dashed line)
933 expected for a well-calibrated model. Panels show instances of (a) under-prediction, and (b)
934 over-prediction of the statistic by the model; (c) under-dispersion, and (d) over-dispersion of
935 statistic in the model simulations.

936

937 Figure 2. (a) Ridgecrest sequence data beginning one week preceding the Mw 6.4 foreshock
938 through the eleven-week evaluation period. Vertical gray dashed lines indicate the starting times
939 of the forecasts. Brown data denote target (test) earthquakes. The forecasts are conditioned on all
940 events until the start time of the forecast. The inset shows the Helmstetter et al. (2006)
941 magnitude-completeness model for the first three days following the Mw 7.1 mainshock. (b)
942 Distribution of spatial seismicity from ComCat during the period shown in (a). The circle shows
943 the spatial region used for the evaluations based on an average Mw 7.1 fault length from Wells
944 and Coppersmith (1994) with a radius of approximately 143 km.

945

946 Figure 3. Synthetic catalog realizations showing 7 days of aftershocks following the M_w 7.1
947 mainshock. (a) ‘Typical’ U3ETAS synthetic catalog, defined as the catalog whose event count
948 lies along the median amongst all simulated catalogs. (b) ‘Extreme’ U3ETAS synthetic catalog,
949 which is defined as the catalog whose event count falls in the uppermost 0.1 percentile of the
950 forecasted number distribution. Notice the triggered ruptures on the Garlock and San Andreas
951 faults that in turn generate aftershocks along these faults. (c) ‘Typical’ synthetic catalog
952 generated by NoFaults and (d) an ‘extreme’ catalog from NoFaults, which lacks triggering of

953 ruptures on prescribed faults resulting in a nearly isotropic aftershock distribution. The ‘extreme’
954 catalogs highlight the predominant differences between these two models and suggest that
955 differences will be most noticeable when large aftershocks occur on mapped faults in U3ETAS.

956

957 Figure 4. Forecasted number distributions and observed cumulative number over the eleven-
958 week evaluation period. The forecasted event count distributions are offset by the number of
959 observed events at the start of the forecast. Forecasted number distributions are plotted at the end
960 of each evaluation period. The vertical extent of the lines indicates the 95-percentile range of the
961 forecasted number distribution. The ‘x’ indicates evaluation periods with observed event counts
962 that fall outside the 95-percentile range of the forecast. (a) Both observed and forecasted catalogs
963 are filtered to threshold magnitudes $M_t(t) = \max(2.5, M_c(t))$ and (b) catalogs are filtered to
964 $M_t(t) = \max(3.5, M_c(t))$. During the first seven-day forecast period, the 95th percentile of the
965 forecasted number distribution for M2.5+ events are 2,482 and 3,906 events for U3ETAS and
966 NoFaults, respectively.

967

968 Figure 5. Aggregate number test results for $M_t(t) = \max(2.5, M_c(t))$ and $M_t(t) =$
969 $\max(3.5, M_c(t))$ magnitude thresholds for U3ETAS and NoFaults for eleven weekly evaluation
970 intervals following the M_w 7.1 mainshock. (a) Quantile scores δ_1 (top) and δ_2 (bottom) for
971 individual weekly evaluation periods. (b) Quantile-quantile plot showing calibration of rate
972 forecasts by comparing quantile scores, γ_N against standard uniform quantiles. The dashed lines
973 indicate 95 percent confidence intervals around the standard uniform quantiles. Thus, U3ETAS
974 and NoFaults overpredict the number of M2.5+ and M3.5+ events during this aftershock
975 sequence.

976

977 Figure 6. (a) Magnitude frequency distribution in $\Delta M = 0.1$ bins aggregated over entire the
978 eleven-week evaluation period. The thin lines approximate the 95% percentile range of the event
979 counts in each magnitude bin. The U3ETAS magnitude frequency distribution shows anti-
980 characteristic behavior through the lack of $M_{6.5+}$ earthquakes as compared with NoFaults. (b)
981 Bin-wise magnitude test statistic aggregated over the entire evaluation period. The circles depict
982 the kernel of d_{obs} for both U3ETAS and NoFaults to show bin-wise contributions to d_{obs} . We
983 find negligible differences between the two models. The solid lines show percentiles from the
984 bin-wise value distribution, for both models.

985

986 Figure 7. Magnitude test results for events with $M_t(t) = (2.5, M_c(t))$ over the full eleven-week
987 evaluation period. (a) Quantile scores are shown for individual week-long evaluation periods.
988 Gray patch depicts the 0.05 significance level for the magnitude test. The largest differences
989 between U3ETAS and NoFaults exist during the first week and become negligible over the
990 remainder of the evaluation period. (b) Calibration of magnitude forecasts by comparing
991 magnitude test quantile scores against standard uniform quantiles. The dashed lines depict 95
992 percent confidence intervals around the standard uniform quantiles.

993

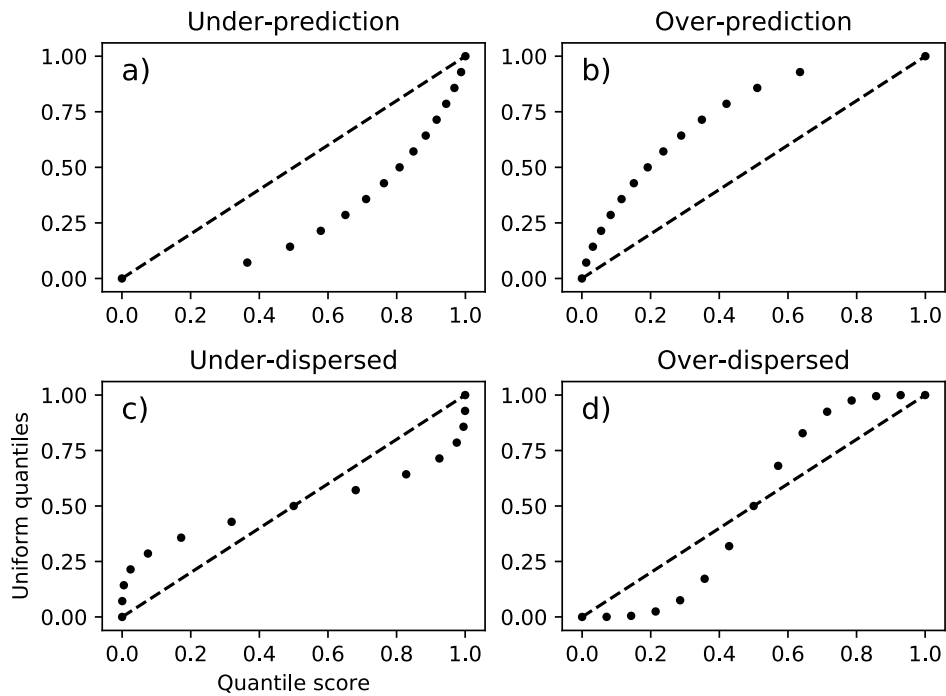
994 Figure 8. Logarithm of the expected event counts per spatial bin per week for U3ETAS (a) and
995 NoFaults (b) for the week-long forecast following the M_w 7.1. The relatively high expected
996 counts along the faults in U3ETAS are controlled by scenarios whose aftershock sequences
997 contain suprasedismogenic ruptures along these faults. In both plots, target events during this
998 period are shown as white circles. The color scale is manually saturated for comparison

999 purposes. The spatial bin with highest rate expects 64.24 and 65.76 events for U3ETAS and
1000 NoFaults, respectively. (c) Evaluation result for the spatial test for U3ETAS (top) and NoFaults
1001 (bottom) for the first evaluation period at seven days after the M_w 7.1 mainshock. $\hat{S}^{(95)}$ denotes
1002 the 95th percentile range of the test distribution of the spatial test statistic, \hat{S}_{obs} is the observed
1003 statistic, and γ_S is the quantile score. (d) Same as (c) except for the pseudo-likelihood test
1004 statistics.

1005
1006 Figure 9. Spatial test and pseudo-likelihood results for events with $M_t(t) = \max(2.5, M_c(t))$
1007 over the complete eleven-week evaluation period. The spatial test and likelihood tests show the
1008 greatest differences between U3ETAS and NoFaults. (a) Quantile scores shown for individual
1009 week-long evaluation periods. The patch depicts the 0.05 significance level for the spatial test.
1010 (b) Calibration of spatial forecasts by comparing quantile scores against standard uniform
1011 quantiles. The dashed lines depict 95 percent confidence intervals around the standard uniform
1012 quantiles.

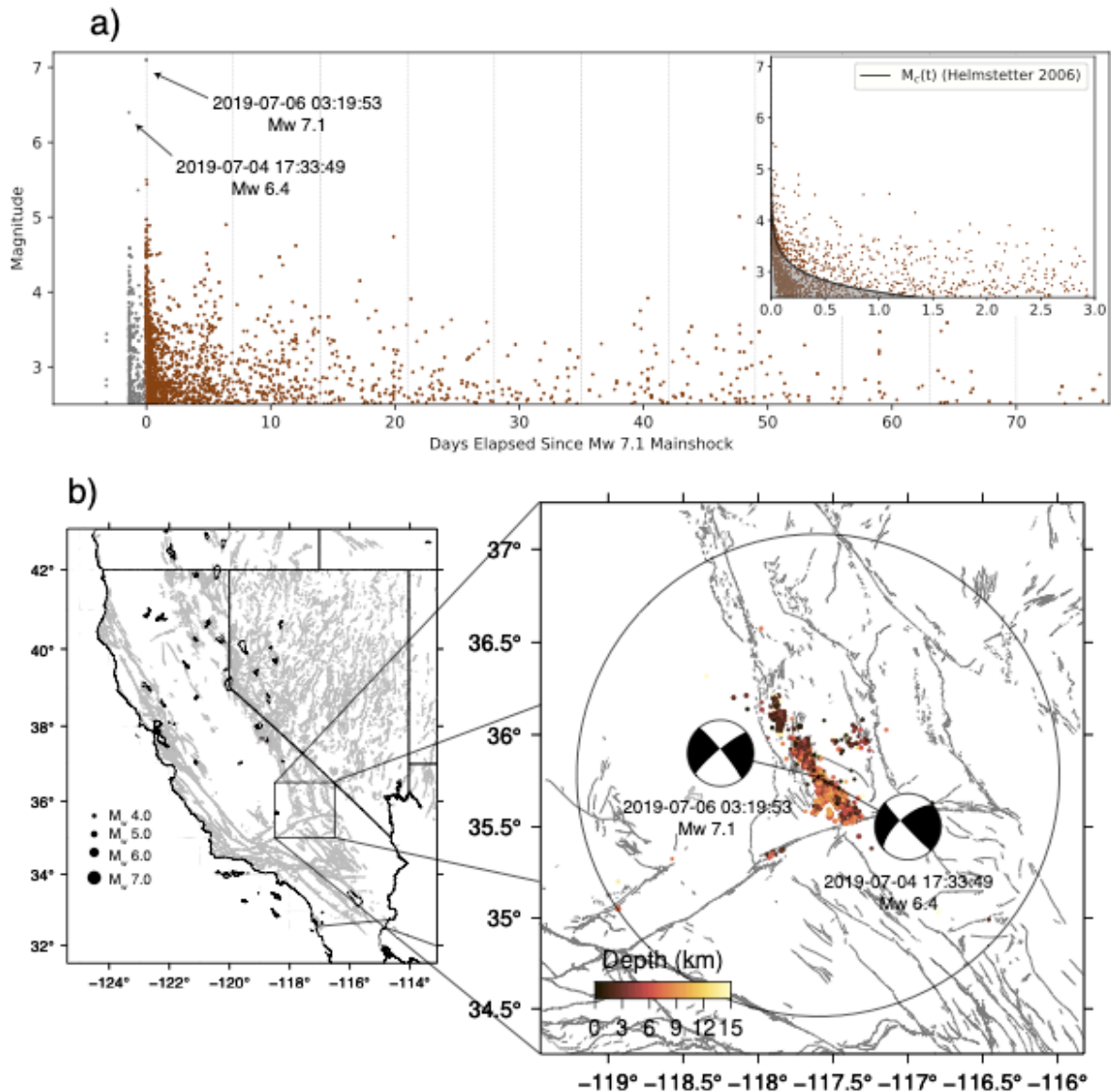
1013
1014 Figure 10. Map of cell-wise spatial pseudo log-likelihood ratios between U3ETAS and NoFaults
1015 for individual evaluation periods ending on (a) day 35, (b) day 49, (c) day 56, and (d) day 63
1016 following the M_w 7.1 mainshock. Maps show the higher rates along faults in U3ETAS.
1017 Evaluation periods at (b) 49 days and (d) 63 days show the largest differences in the observed
1018 spatial statistic, which is calculated only from spatial cells where events occur, while periods
1019 ending on days 35 and 56 show a negligible difference in the spatial statistic. This highlights
1020 how spatial test results are sensitive to events occurring on modeled U3ETAS faults and that
1021 such events are required to discern between the models. The color

1022 scale is manually saturated between -0.05 and 0.05 to help comparisons; and dots show
1023 locations of target events
1024

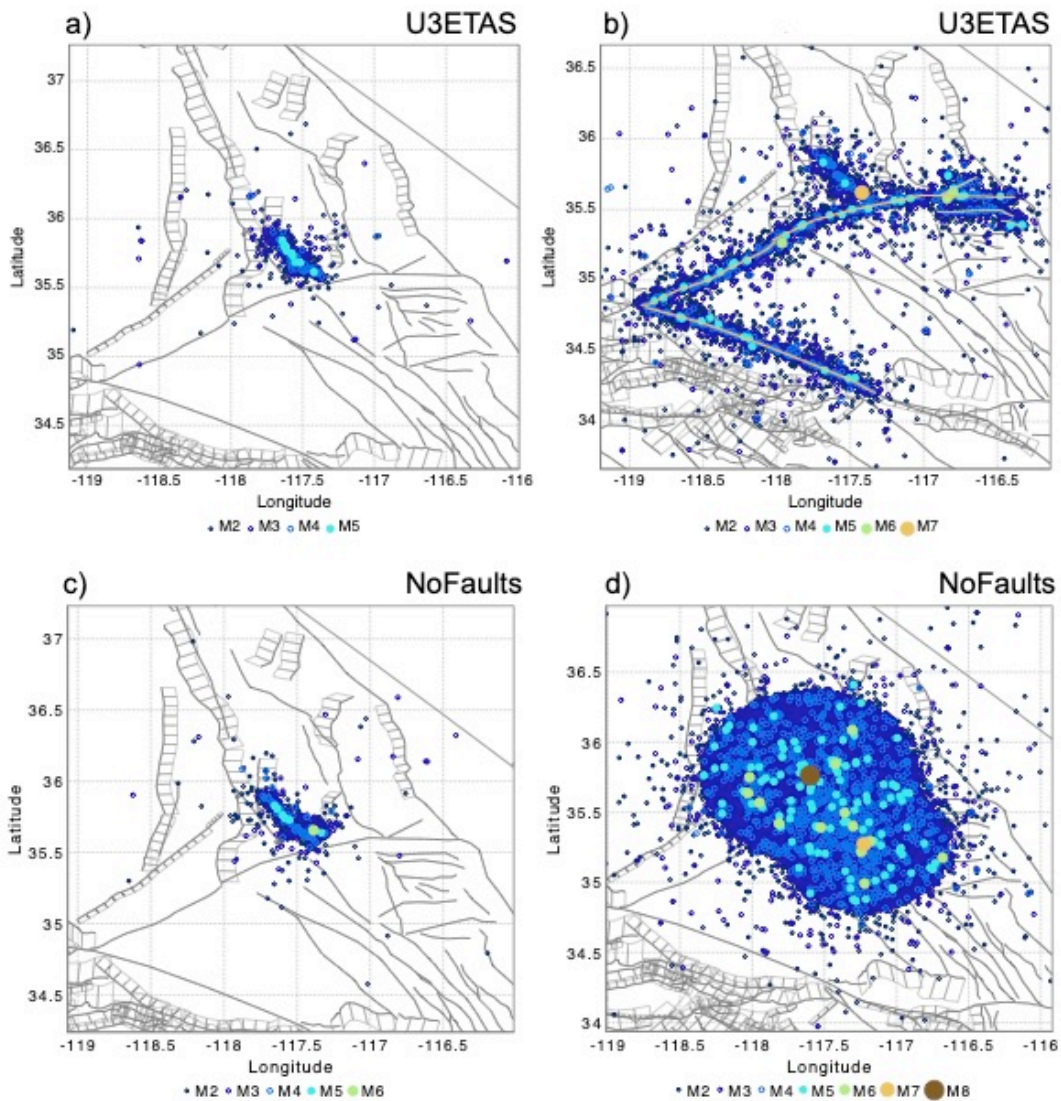


1026

1027 Figure 1. Schematic of cumulative distribution of quantile scores for a test statistic calculated over multiple test
 1028 periods (points) as compared with the ideal uniform distribution (dashed line) expected for a well-calibrated model.
 1029 Panels show instances of (a) under-prediction, and (b) over-prediction of the statistic by the model; (c) under-
 1030 dispersion, and (d) over-dispersion of statistic in the model simulations.



1031
 1032 Figure 2. (a) Ridgecrest sequence data beginning one week preceding the M_w 6.4 foreshock through the eleven-week
 1033 evaluation period. Vertical gray dashed lines indicate the starting times of the forecasts. Brown data denote target
 1034 (test) earthquakes. The forecasts are conditioned on all events until the start time of the forecast. The inset shows the
 1035 Helmstetter et al. (2006) magnitude-completeness model for the first three days following the M_w 7.1 mainshock. (b)
 1036 Distribution of spatial seismicity from ComCat during the period shown in (a). The circle shows the spatial region
 1037 used for the evaluations based on an average M_w 7.1 fault length from Wells and Coppersmith (1994) with a radius
 1038 of approximately 143 km.



1039

1040

1041

1042

1043

1044

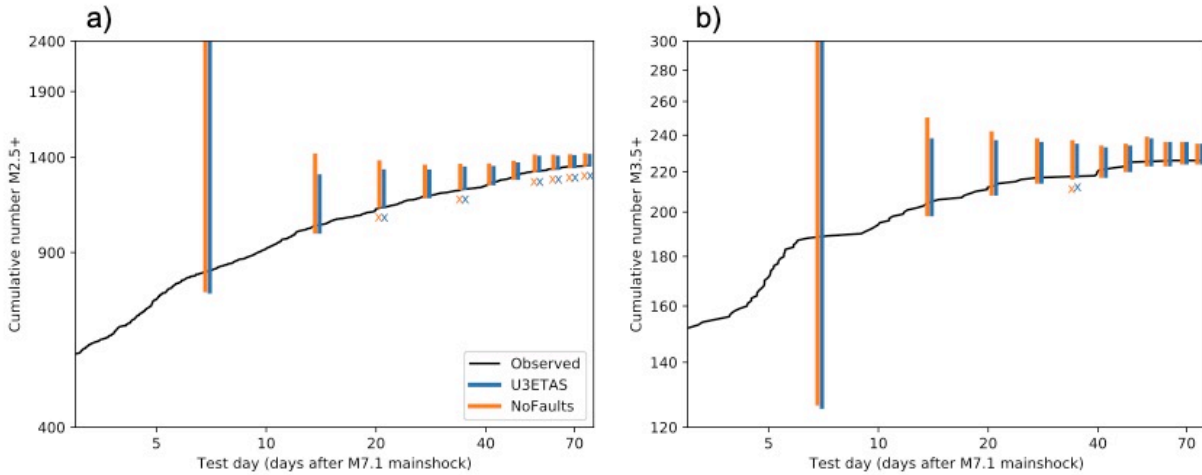
1045

1046

1047

1048

Figure 3. Synthetic catalog realizations showing 7 days of aftershocks following the M_w 7.1 mainshock. (a) ‘Typical’ U3ETAS synthetic catalog, defined as the catalog whose event count lies along the median amongst all simulated catalogs. (b) ‘Extreme’ U3ETAS synthetic catalog, which is defined as the catalog whose event count falls in the uppermost 0.1 percentile of the forecasted number distribution. Notice the triggered ruptures on the Garlock and San Andreas faults that in turn generate aftershocks along these faults. (c) ‘Typical’ synthetic catalog generated by NoFaults and (d) an ‘extreme’ catalog from NoFaults, which lacks triggering of ruptures on prescribed faults resulting in a nearly isotropic aftershock distribution. The ‘extreme’ catalogs highlight the predominant differences between these two models and suggest that differences will be most noticeable when large aftershocks occur on mapped faults in U3ETAS.



1049

1050 Figure 4. Forecasted number distributions and observed cumulative number over the eleven-week evaluation period.

1051 The forecasted event count distributions are offset by the number of observed events at the start of the forecast.

1052 Forecasted number distributions are plotted at the end of each evaluation period. The vertical extent of the lines

1053 indicates the 95-percentile range of the forecasted number distribution. The 'x' indicates evaluation periods with

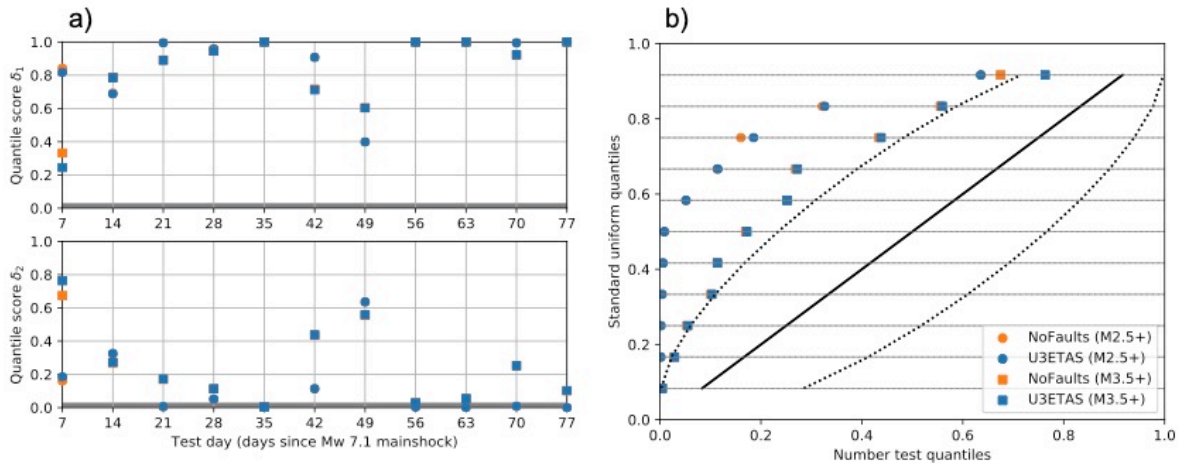
1054 observed event counts that fall outside the 95-percentile range of the forecast. (a) Both observed and forecasted

1055 catalogs are filtered to threshold magnitudes $M_t(t) = \max(2.5, M_c(t))$ and (b) catalogs are filtered to

1056 $M_t(t) = \max(3.5, M_c(t))$. During the first seven-day forecast period, the 95th percentile of the forecasted number

1057 distribution for $M_{2.5+}$ events are 2,482 and 3,906 events for U3ETAS and NoFaults, respectively.

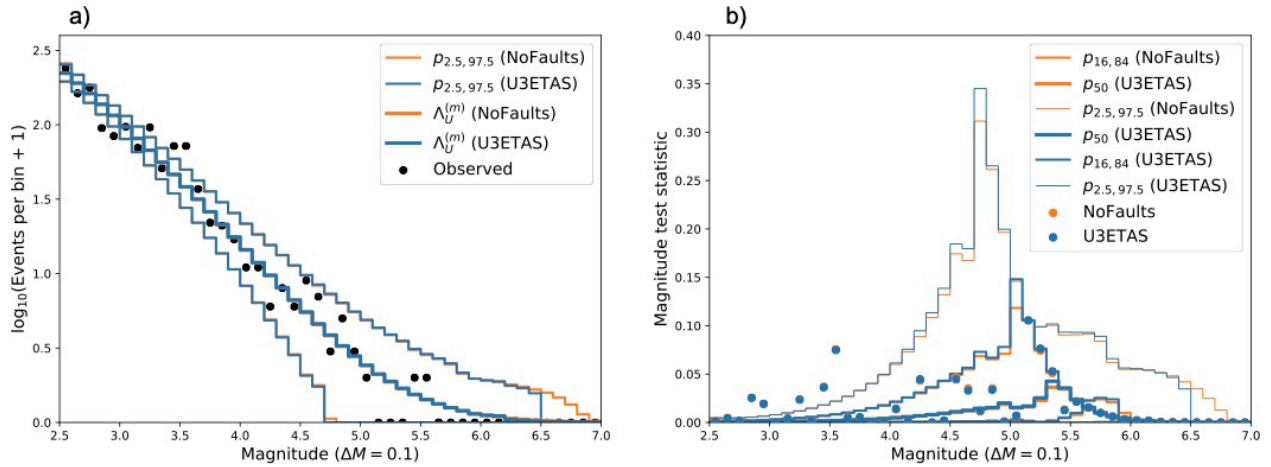
1058



1059

1060 Figure 5. Aggregate number test results for $M_t(t) = \max(2.5, M_c(t))$ and $M_t(t) = \max(3.5, M_c(t))$ magnitude
 1061 thresholds for U3ETAS and NoFaults for eleven weekly evaluation intervals following the M_w 7.1 mainshock. (a)
 1062 Quantile scores δ_1 (top) and δ_2 (bottom) for individual weekly evaluation periods. (b) Quantile-quantile plot
 1063 showing calibration of rate forecasts by comparing quantile scores, γ_N against standard uniform quantiles. The
 1064 dashed lines indicate 95 percent confidence intervals around the standard uniform quantiles. Thus, U3ETAS and
 1065 NoFaults overpredict the number of M2.5+ and M3.5+ events during this aftershock sequence.

1066



1067

1068 Figure 6. (a) Magnitude frequency distribution in $\Delta M = 0.1$ bins aggregated over entire the eleven-week evaluation

1069 period. The thin lines approximate the 95% percentile range of the event counts in each magnitude bin. The

1070 U3ETAS magnitude frequency distribution shows anti-characteristic behavior through the lack of M6.5+

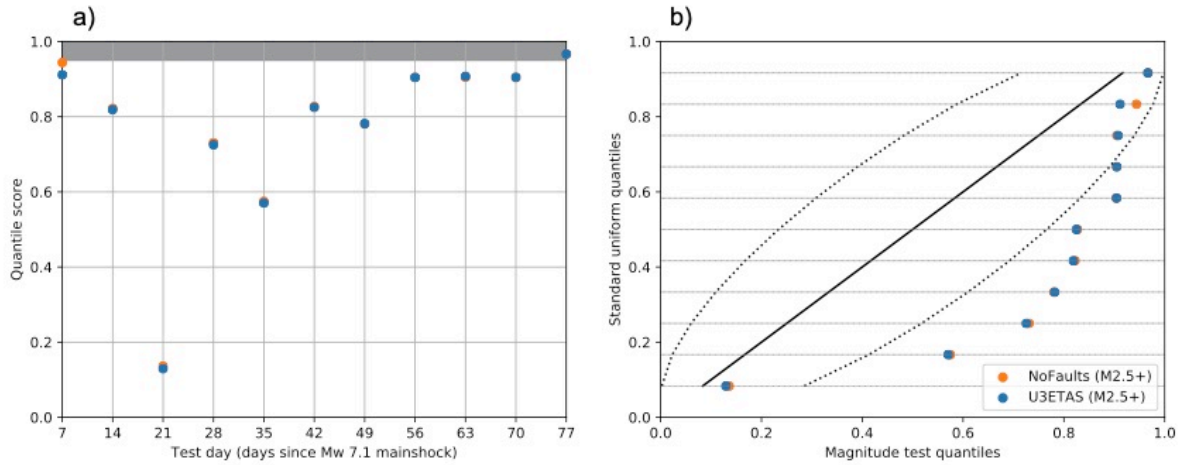
1071 earthquakes as compared with NoFaults. (b) Bin-wise magnitude test statistic aggregated over the entire evaluation

1072 period. The circles depict the kernel of d_{obs} for both U3ETAS and NoFaults to show bin-wise contributions to d_{obs} .

1073 We find negligible differences between the two models. The solid lines show percentiles from the bin-wise value

1074 distribution, for both models.

1075



1076

1077 Figure 7. Magnitude test results for events with $M_t(t) = (2.5, M_c(t))$ over the full eleven-week evaluation period.

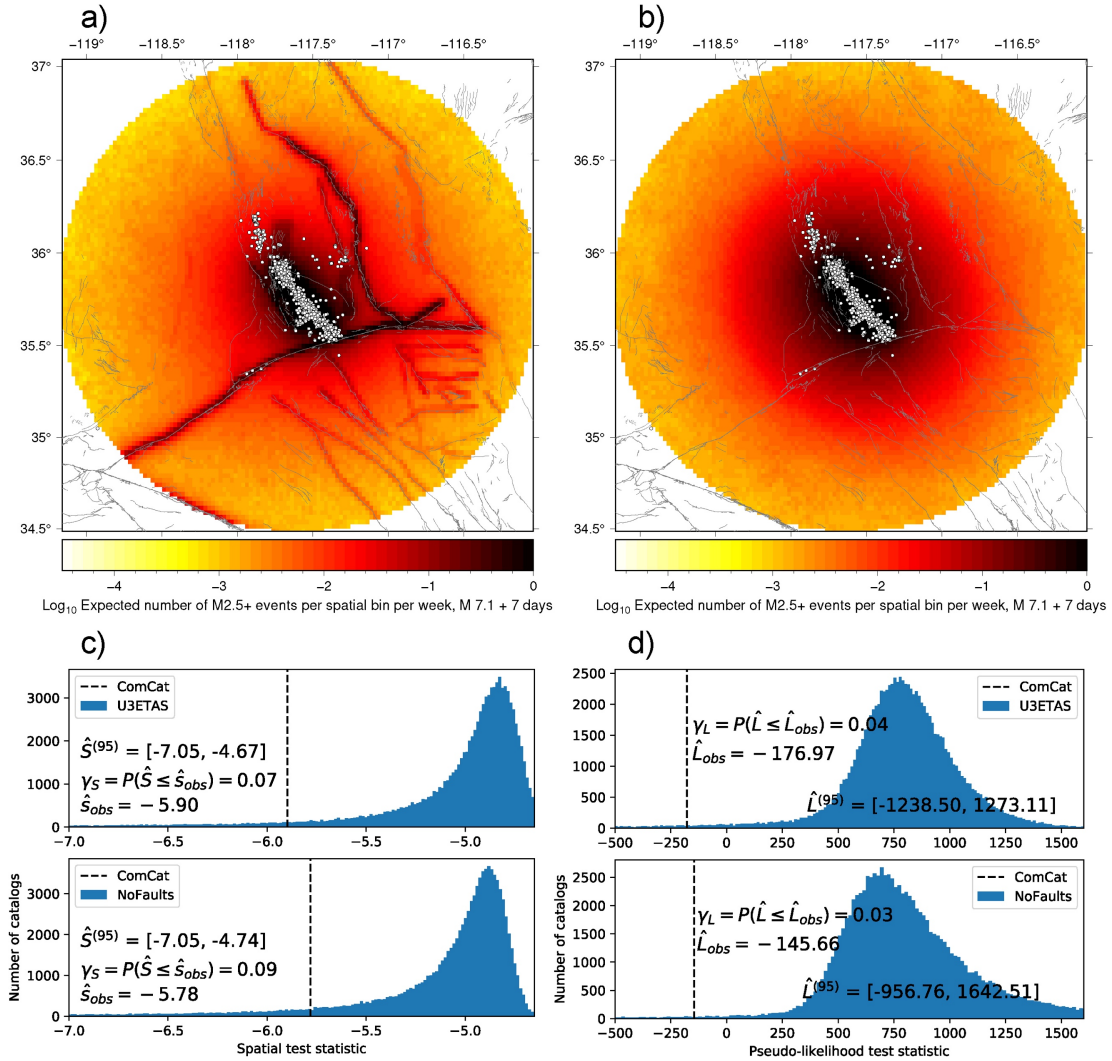
1078 (a) Quantile scores are shown for individual week-long evaluation periods. Gray patch depicts the 0.05 significance

1079 level for the magnitude test. The largest differences between U3ETAS and NoFaults exist during the first week and

1080 become negligible over the remainder of the evaluation period. (b) Calibration of magnitude forecasts by comparing

1081 magnitude test quantile scores against standard uniform quantiles. The dashed lines depict 95 percent confidence

1082 intervals around the standard uniform quantiles.



1083

1084

1085

1086

1087

1088

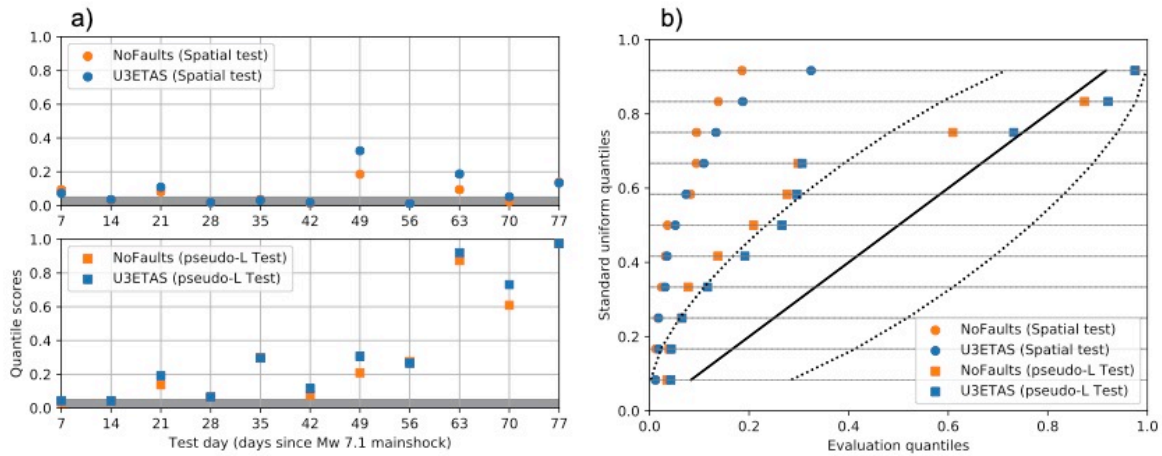
1089

1090

1091

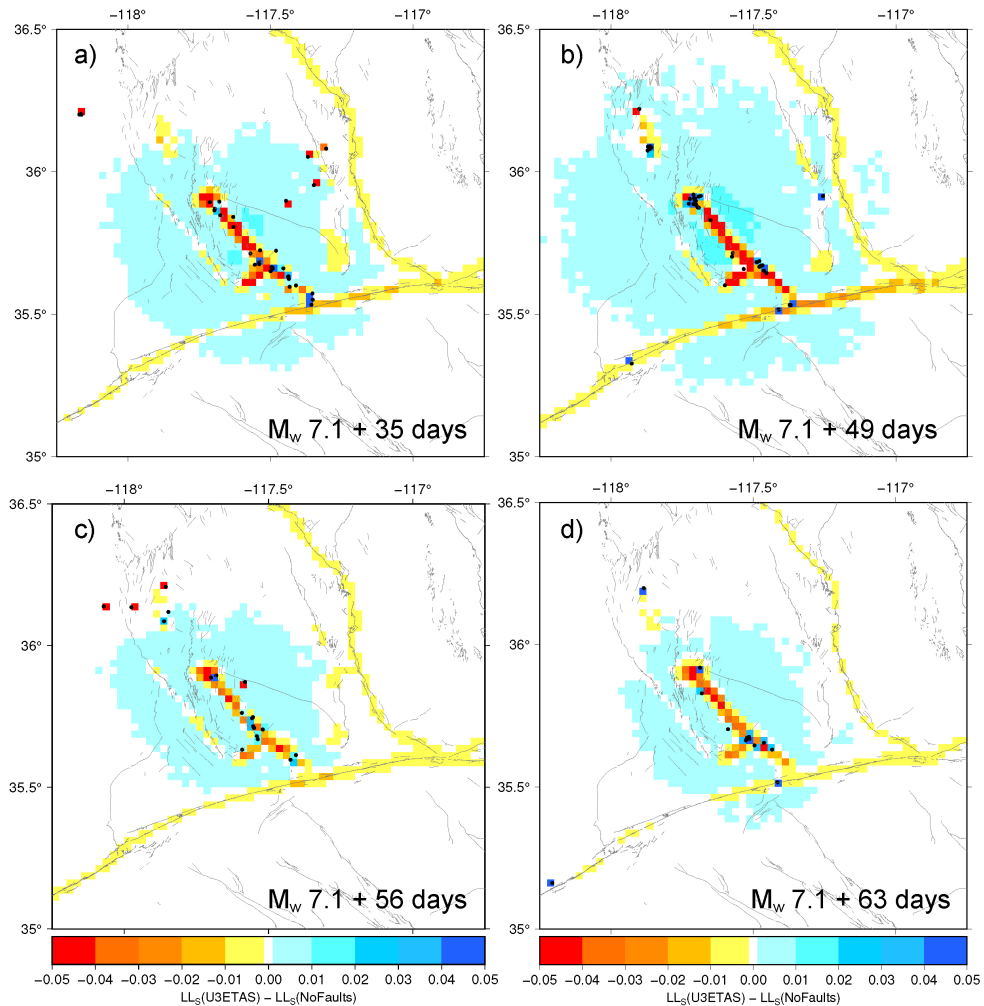
1092

Figure 8. Logarithm of the expected event counts per spatial bin per week for U3ETAS (a) and NoFaults (b) for the week-long forecast following the M_w 7.1. The relatively high expected counts along the faults in U3ETAS are controlled by scenarios whose aftershock sequences contain suprasedismogenic ruptures along these faults. In both plots, target events during this period are shown as white circles. The color scale is manually saturated for comparison purposes. The spatial bin with highest rate expects 64.24 and 65.76 events for U3ETAS and NoFaults, respectively. (c) Evaluation result for the spatial test for U3ETAS (top) and NoFaults (bottom) for the first evaluation period at seven days after the M_w 7.1 mainshock. $\hat{S}^{(95)}$ denotes the 95th percentile range of the test distribution of the spatial test statistic, \hat{s}_{obs} is the observed statistic, and γ_S is the quantile score. (d) Same as (c) except for the pseudo-likelihood test statistics.



1093

1094 Figure 9. Spatial test and pseudo-likelihood results for events with $M_t(t) = \max(2.5, M_c(t))$ over the complete
 1095 eleven-week evaluation period. The spatial test and likelihood tests show the greatest differences between U3ETAS
 1096 and NoFaults. (a) Quantile scores shown for individual week-long evaluation periods. The patch depicts the 0.05
 1097 significance level for the spatial test. (b) Calibration of spatial forecasts by comparing quantile scores against
 1098 standard uniform quantiles. The dashed lines depict 95 percent confidence intervals around the standard uniform
 1099 quantiles.



1100

1101 Figure 10. Map of cell-wise spatial pseudo log-likelihood ratios between U3ETAS and NoFaults for individual
 1102 evaluation periods ending on (a) day 35, (b) day 49, (c) day 56, and (d) day 63 following the M_w 7.1 mainshock.

1103 Maps show the higher rates along faults in U3ETAS. Evaluation periods at (b) 49 days and (d) 63 days show the
 1104 largest differences in the observed spatial statistic, which is calculated only from spatial cells where events occur,

1105 while periods ending on days 35 and 56 show a negligible difference in the spatial statistic. This highlights how

1106 spatial test results are sensitive to events occurring on modeled U3ETAS faults and that such events are required to

1107 discern between the models. The color scale is manually saturated between -0.05 and 0.05 to help comparisons; and

1108 dots show locations of target events.

1109

1110 **Author Contact Information**

1111

1112 Maximilian J. Werner, School of Earth Sciences, University of Bristol, Wills Memorial Building,
1113 Queens Road, Bristol BS8 1RJ, United Kingdom

1114

1115 Warner Marzocchi, Università degli Studi di Napoli Federico II, Corso Umberto I, 40, 80138
1116 Napoli NA, Italy

1117

1118 David Rhoades, GNS Science, 1 Fairway Drive, Avalon, Lower Hutt 5011, New Zealand

1119

1120 David Jackson, Earth, Planetary, and Space Sciences University of California, Los Angeles 595
1121 Charles Young Drive East Box 951567 Los Angeles, CA 90095-1567

1122

1123 Kevin Milner, University of Southern California, Southern California Earthquake Center, 3651
1124 Trousdale Parkway, Los Angeles, CA 90089

1125

1126 Ned Field, United States Geological Survey, 1711 Illinois St Golden, CO 80401

1127

1128 Andrew Michael, United States Geological Survey, U.S. Geological Survey, 350 N. Akron Road
1129 Moffett Field, CA 94035

1130