

1 **pyCSEP: A Python Toolkit for Earthquake**

2 **Forecast Developers**

3 William H. Savran¹, José A. Bayona², Pablo Iturrieta³, Khawaja M. Asim³, Han
4 Bao⁴, Kirsty Bayliss⁵, Marcus Herrmann⁶, Danijel Schorlemmer³, Philip J.
5 Maechling¹, and Maximilian J. Werner²

6 ¹University of Southern California

7 ²University of Bristol

8 ³GFZ Potsdam

9 ⁴University of California, Los Angeles

10 ⁵University of Edinburgh

11 ⁶University of Naples 'Federico II', Italy

12 **Corresponding Author**

13 William H. Savran, Zumberge Hall of Science (ZHS), 3651 Trousdale Pkwy, Los Angeles, CA 90089-0740,
14 wsavran@usc.edu

15 **Declaration of Competing Interests**

16 The authors acknowledge there are no conflicts of interest recorded.

17 **ABSTRACT**

18 The Collaboratory for the Study of Earthquake Predictability (CSEP) is an open and global community
19 whose mission is to accelerate earthquake predictability research through rigorous testing of probabilistic
20 earthquake forecast models and prediction algorithms. pyCSEP supports this mission by providing open-
21 source implementations of useful tools for evaluating earthquake forecasts. pyCSEP is a Python package
22 that contains the following modules: (1) earthquake catalog access and processing, (2) representations
23 of probabilistic earthquake forecasts, (3) statistical tests for evaluating earthquake forecasts, and (4)
24 visualization routines and various other utilities. Most significantly, pyCSEP contains several statistical
25 tests needed to evaluate earthquake forecasts, which can be forecasts expressed as expected earthquake
26 rates in space-magnitude bins or specified as large sets of simulated catalogs (which includes candidate
27 models for governmental operational earthquake forecasting). To showcase how pyCSEP can be used to
28 evaluate earthquake forecasts, we have provided a reproducibility package that contains all the components
29 required to recreate the figures published in this article. We recommend that interested readers work
30 through the reproducibility package alongside this manuscript. By providing useful tools to earthquake
31 forecast modelers and facilitating an open-source software community, we hope to broaden the impact
32 of the Collaboratory for the Study of Earthquake Predictability (CSEP) and further promote earthquake
33 forecasting research.

34 **INTRODUCTION**

35 **The Collaboratory for the Study of Earthquake Predictability**

36 CSEP emerged from the need to place the research field on more robust methodological footing to help
37 overcome the negative sentiment surrounding earthquake prediction efforts (e.g., Geller, 1997). CSEP
38 formed as a collaboration to assess earthquake predictability and provide users of earthquake forecasts
39 with confidence about forecast skill and performance (e.g., government agencies that issue operational
40 earthquake forecasts; Jordan and Jones, 2010; Jordan et al., 2011; Marzocchi et al., 2014). Past efforts
41 were stymied by a range of problems that resulted in a lack of both reproducibility (the inability to
42 regenerate previously issued forecasts, predictions, or test results) and replicability (the inability to reach
43 the same conclusion about a model's predictive skill from different data; Stodden et al., 2018; National
44 Academies of Sciences, Engineering, and Medicine and others, 2019). The peer-review process was
45 frequently insufficient to ensure these necessary standards, an experience mirrored in other empirical
46 research fields (Baker, 2016). Meaningful prospective evaluations require sufficient data, which may take
47 several decades or more to collect in certain regions, especially for large earthquakes. CSEP's multi-region
48 approach and global experiments, which would not be possible without its international collaboration,

49 help alleviate this limitation (e.g., Bird et al., 2015). Although progress in forecast testing may be limited
50 by time, even a few years of data help scientists to falsify certain hypotheses that are inconsistent with
51 observations (Dekel and Feinberg, 2006).

52 The main pillar of CSEP's approach is the prospective testing of forecasts (i.e., against future
53 observations) in reproducible and transparent forecast experiments carefully designed by the community.
54 Prospective evaluations require that forecasts are unambiguously testable, with all model parameters,
55 forecast specifications, and qualifying target data sources specified in advance, preferably before testing
56 observations were made (Schorlemmer and Gerstenberger, 2007; Schorlemmer et al., 2018). This ensures
57 a zero-degree-of-freedom, independent test of a model's or algorithm's performance.

58 Starting in 2007, CSEP has managed testing centers that autonomously run prospective forecast
59 experiments (Schorlemmer and Gerstenberger, 2007). In these, automated dispatchers run forecast models
60 to generate forecasts and evaluate them against prospective data (Zechar et al., 2010). Testing centers
61 existed in California, New Zealand, Italy, Japan, and China, and together hosted over 400 models and
62 model versions in a variety of tectonic settings and at a global scale (e.g., Field, 2007; Marzocchi et al.,
63 2014; Tsuruoka et al., 2012; Zechar et al., 2013; Taroni et al., 2018; Strader et al., 2018; Rhoades et al.,
64 2018; Eberhard et al., 2012; Bayona et al., 2021). Through this major community effort, CSEP has
65 provided new insights into the predictability of earthquakes, provided independent assessments of the
66 predictive skills of a range of scientific hypotheses of seismogenesis, galvanised model improvements and
67 motivated new research into evaluation methods (Schorlemmer et al., 2018).

68 After a decade of operating the CSEP testing centers, it became apparent that the monolithic software
69 design was too strongly entangled with the system architecture and data bookkeeping to support the new
70 types of forecast experiments that the CSEP community would like to conduct (Schorlemmer et al., 2018).
71 CSEP software has always been open-source and accessible; however, in practice, the code was difficult to
72 use by individual researchers. Specifically, the testing center software coupled the evaluation routines with
73 the system architecture making it difficult to use them outside of the testing center context. We developed
74 pyCSEP as the first of many steps to modernize CSEP testing centers and experiments. Modern testing
75 centers should use pyCSEP as a library, decoupling the testing center architecture from the evaluation
76 routines. Additionally, they should follow modern open-science principles to ensure that experiment
77 results are versioned and openly available to the public (e.g., Wilkinson et al., 2016). Testing centers
78 are crucial for addressing the replicability of experimental results, because long-standing prospective
79 experiments are required to capture the time-scales needed for model improvements and updates.

80 **pyCSEP: A Python Toolkit for Earthquake Forecast Developers**

81 Strengthening the collaborative aspects of the CSEP community and increasing the sustainability of CSEP
82 activities, requires a new and collaborative mode of software development with the goal of a flexible,
83 open-source, and community-based processing toolkit. Building sustainable research software requires a
84 community that bridges software engineers and scientists (Anzt et al., 2021). This open-source approach
85 is ideal for research software, as it allows for transparent, extendable code development by the research
86 community that is using the software. It allows practitioners of the code to implement new features and
87 identify potential issues in the software, and become engaged with the development process creating a
88 net benefit for all involved members. We conceived of the open-source pyCSEP toolkit to address this
89 limitation and to create a software community to promote earthquake forecasting research.

90 At its core, pyCSEP re-implements software running in CSEP testing centers as an open-source Python
91 package, but is already rapidly expanding beyond this. pyCSEP is designed so researchers can evaluate
92 earthquake forecasting models with minimal effort using a beginner-friendly, object-oriented interface.
93 pyCSEP's modular structure allows for easy extensibility (Fig. 1). We encourage researchers to contribute
94 code directly to the toolkit. To enable reproducible research, we strive for collaboratively developed code
95 that is readable, well documented, and, most importantly, vetted. The source code can be found in the
96 GitHub repository for this project (see link in *Data and Resources* section). Savran et al. (2022) provides
97 a brief overview of the motivation for developing pyCSEP. The review process for that article focuses
98 on software development best-practices, and examines the software repository and documentation. This
99 article complements the software focused publication by providing more thorough explanations of the
100 functionality of the software and providing the accompanying reproducibility package.

101 **Software Development Principles**

102 We incorporate several best-practices used by many open-source software projects (e.g., Hunter, 2007;
103 McKinney, 2010; pandas Development Team, 2020) into our development process. In the code repository,
104 we use continuous integration (CI) tools to ensure all new code contributions build successfully and pass
105 unit tests. CI tools trigger workflows in the software repository to run development tasks automatically.
106 The CI tools also build and publish the online documentation (link in *Data and Resources* section).
107 These workflows trigger automatically when changes are made to the 'main' branch of the repository,
108 or when new contributions have been submitted as pull requests in GitHub. We follow the SemVer
109 (<https://semver.org>) guidelines for software versioning. New releases are made available on `PyPI` and
110 `conda-forge` and can be installed using the package managers `pip` or `conda`. Additionally, pyCSEP
111 strives to meet the target best practices as proposed by the Computational Infrastructure for Geodynamics
112 (link in *Data and Resources* section).

113 **Reproducibility of Forecasting Experiments**

114 In CSEP testing centers, experiment components (e.g., model software, input data, forecasts, target data,
115 and test results) were stored on CSEP servers with no external access (Schorlemmer and Gerstenberger,
116 2007). This approach ensured the integrity and reproducibility of the experiments, but required substantial
117 data management and systems administration resources. The controlled environment of CSEP testing
118 centers also made it difficult to share experimental results. The recent proliferation of freely available
119 online data storage and management tools provide an effective alternative for storing experiment data
120 and code. We encourage the use of these tools to create reproducibility packages (Krafczyk et al., 2021)
121 for publications of earthquake forecasting experiments. A reproducibility package contains the software,
122 data, and other experiment artifacts required to exactly reproduce published results. To illustrate this idea
123 and provide an introduction to pyCSEP, we provide an example reproducibility package for this article
124 (link in *Data and Resources* section).

125 **PYCSEP SOFTWARE**

126 pyCSEP provides an open-source implementation of several peer-reviewed statistical tests for evaluating
127 probabilistic earthquake forecasts (Schorlemmer et al., 2007; Zechar et al., 2010; Rhoades et al., 2011;
128 Werner et al., 2011; Savran et al., 2020). The design includes core classes that represent earthquake
129 forecasts, catalogs, and spatial regions (Fig. 1). Higher-level functions using these classes are implemented
130 to provide a simple interface to analyze forecasting models. Overall, the software design is modular to
131 accommodate new forecast representations and evaluation types. Where possible, we integrate popular
132 Python libraries such as numpy (Harris et al., 2020), matplotlib (Hunter, 2007), and pandas (pandas
133 Development Team, 2020; McKinney, 2010) to allow users to easily include pyCSEP in existing scripts
134 and workflows. pyCSEP also contains routines for working with and visualizing earthquake forecasts and
135 catalogs. Also, general users of earthquake catalogs and gridded data sets may find useful utilities in the
136 package.

137 **Getting started with pyCSEP**

138 The most straightforward way to install pyCSEP is using the `conda` package manager, and installing
139 the most recent release from `conda-forge`. Users can obtain `conda` through the Anaconda or
140 `miniconda` distributions. pyCSEP issues regular releases to PyPI and `conda-forge`. The latest
141 release can be installed using:

```
142     conda install --channel conda-forge pycsep
```

143 The online documentation provides detailed installation instructions and examples that assist new

144 users through tasks such as evaluating grid- and catalog-based earthquake forecasts, working with catalogs,
145 and various plotting tasks (link in *Data and Resources* section).

146 **Core Classes**

147 The following subsections present a more technical introduction to the core classes in pyCSEP (see Fig.
148 1). Fig. 1 indicates how important methods and classes are related in the code, and can be used as a
149 reference. We recommend interested readers to get started with pyCSEP by following the examples in
150 the online documentation, and working through the reproducibility package (see Section *Reproducibility*
151 *Packages* for this manuscript; link in *Data and Resources* section).

152 **Regions**

153 Regions are used to define the spatial cells of an earthquake forecast. In practice, they are used to bin, or
154 discretize, an earthquake catalog into these spatial cells. Regions are fundamental in defining earthquake
155 forecasts and preparing observed catalogs for evaluation (Fig. 1). In practice, a region represents a
156 mapping between a list of spatial cells and spatial points. This mapping associates each point with its
157 corresponding cell in the spatial region. There is a many-to-one relationship between points and spatial
158 cells. Each point can only be associated with a single cell; however, a cell can contain many points.
159 pyCSEP defines a standard region on a regular Cartesian grid whose cells have dimensions of $0.1^\circ \times 0.1^\circ$
160 in latitude and longitude. However, the dimensions of the cells are configurable within pyCSEP. To allow
161 for easy interoperability with previous experiments, pyCSEP currently provides predefined regions for
162 California (Fig. 2a), Italy (Fig. 2b) and the global testing region (not pictured). These can be accessed via
163 simple function calls (e.g., `california_relm_region()`, see Fig. 1).

164 pyCSEP provides a class named `CartesianGrid2D` to represent the standard region used in CSEP
165 experiments (e.g., Schorlemmer and Gerstenberger, 2007; Taroni et al., 2018). `CartesianGrid2D`
166 implements the mapping so points can be correctly associated with the corresponding Cartesian spatial
167 cells. The class provides flexibility for creating different regions by supplying a list containing the
168 lower-left origin of each cell by calling the `from_origins()` class method. The cells are defined
169 such that the lower and left-most edges are inclusive. Functionality for non-regular grids is not currently
170 implemented in the toolkit; however, the object-oriented implementation of the region class allows for
171 non-regular grids to be easily accommodated in the future.

172 Magnitude ranges are defined using a list containing the left bin edges, and require no additional
173 classes. The magnitude bin edges should be made accessible through the magnitude member of the
174 region classes. The `regions.create_space_magnitude_region` function provides a method
175 to associate a discretized magnitude range with a particular spatial region.

176 **Forecasts**

177 Currently, pyCSEP supports two types of probabilistic earthquake forecasts (see Fig. 1). First, we
178 support grid-based forecasts that express expected rates of earthquakes within discrete space-time-
179 magnitude bins (e.g., Schorlemmer and Gerstenberger, 2007). A grid-based forecast is defined by the
180 *GriddedForecast* class. This class is composed of two main data attributes: 1) a 2D numpy array
181 that stores expected rates in space-magnitude bins, and 2) a pyCSEP region class that defines the space-
182 magnitude cells of the forecast. Standard CSEP gridded forecasts use the *CartesianGrid2D* to define
183 this mapping. Each forecast is considered to span a discrete time period, where the expected rate is
184 based on this period. Thus, time-dependent forecasts with multiple periods require individual instances of
185 the *GriddedForecast* class. Additional methods such as *target_event_rates()* are provided
186 by the forecast class, and allow the users to retrieve event rates as defined by the forecast. Grid-based
187 forecasts can be loaded from disk using the *load_gridded_forecast()* function defined in the
188 top-level package.

189 The *CatalogForecast* class defines the second supported forecast type: catalog-based forecasts.
190 This class represents forecasts that are defined by a list of earthquake catalogs (e.g., *CSEPCatalog*
191 or *UCERF3Catalog*) and a region (e.g., *CartesianGrid2D*). The class provides an iterator imple-
192 menting a user-defined set of catalog filters that apply automatically to each catalog in the forecast. Also,
193 this implementation allows for working with large UCERF3-ETAS (or other) forecasts by loading the
194 catalogs on demand. This is known as ‘lazy’ loading. *CatalogForecast* objects can be loaded from
195 disk using the *load_catalog_forecast()* function defined in the top-level package. Fig. 3 shows
196 an example of an UCERF3-ETAS forecast made during the 2019 Ridgecrest sequence. The reader will
197 find examples of working with grid-based and catalog-based forecasts in the *Tutorials* section of the
198 online documentation and the reproducibility package.

199 **Evaluations**

200 CSEP has lead research efforts into developing forecast evaluation methods, tests, and performance
201 measures of probabilistic earthquake forecasts (e.g., Schorlemmer et al., 2007; Werner and Sornette, 2008;
202 Zechar et al., 2010; Zechar and Jordan, 2010; Zechar and Zhuang, 2010; Rhoades et al., 2011; Werner
203 et al., 2011; Marzocchi et al., 2012; Schneider et al., 2014; Gordon et al., 2015; Molchan et al., 2017;
204 Savran et al., 2020, and many others). Different tests are used to address various hypotheses underlying
205 the forecasts they are evaluating. pyCSEP currently contains a selection of consistency tests (comparing
206 forecasts with data) and comparative tests (comparing models against each other on the basis of the data)
207 for both grid-based forecasts and catalog-based forecasts. Different forecast formats require different
208 evaluation methods. Grid-based forecasts use a set of evaluations that based on the Poisson likelihood

209 function (Schorlemmer et al., 2007; Zechar et al., 2010), whereas catalog-based forecasts build empirical
210 distributions to sample the uncertainty contained within the forecast (Nandan et al., 2019; Savran et al.,
211 2020). The Poisson assumption has been widely criticized (Lombardi and Marzocchi, 2010a; Werner and
212 Sornette, 2008) and pyCSEP was designed to accommodate evaluation with different likelihood functions
213 (Bayona et al., 2022). We explain the evaluation methods implemented in pyCSEP below. Evaluations for
214 grid-based forecasts are implemented in the module *poisson_evaluations*, and for catalog-based
215 forecasts in the *catalog_evaluations* module (Fig. 1). Examples on how to evaluate grid- and
216 catalog-based forecasts are shown in the *Tutorial* section of the online documentation. We provide an
217 in-depth explanation of the evaluations along with working code examples in the Electronic Supplement
218 to this article.

219 For grid-based forecasts, CSEP tests assess the consistency between the observed and the expected
220 number, spatial, magnitude, and likelihood distributions of earthquakes, assuming that seismicity in
221 space-magnitude cells is independent and Poisson-distributed (Zechar et al., 2010; Werner et al., 2011;
222 Rhoades et al., 2011). In the following paragraphs, we provide a high-level overview of the test methods
223 available for grid-based forecasts followed by a brief description of the consistency tests for catalog-based
224 forecasts.

225 **Number test** The number (N) test (Schorlemmer et al., 2007; Zechar et al., 2010) evaluates if the
226 total number of observed earthquakes (N_{obs}) falls within the 95% predictive distribution of the forecast
227 distribution, with the expected rate, N_{fore} , equal to the sum of forecasted rates in each space-magnitude bin.
228 Fig. 4 shows the N-test result for time-independent forecasts from the Regional Earthquake Likelihood
229 Model (RELM) experiment that were originally published by Zechar et al. (2013).

230 **Spatial test** The spatial (S) test (Zechar et al., 2010) evaluates how well a forecast explains the spatial
231 distribution of earthquakes. One first sums the expected rates in each spatial cell over the magnitude
232 bins to isolate the spatial component of the forecast, and normalizes the resulting spatial rates to the
233 total number of target observations. Next, one computes the (spatial) joint log-likelihood in each cell by
234 evaluating the Poisson likelihood function in each cell, and summing the spatial log-likelihoods over the
235 entire testing region. To assess whether this observed log-likelihood score could have been generated
236 by the forecast, we obtain the distribution of spatial log-likelihood scores consistent with the forecast
237 through simulation. In this and the following two tests, the number of simulated earthquakes is fixed to
238 N_{obs} to remove the dependency on the forecasted rate. To assess the consistency between the observed
239 locations and the spatial forecast, we examine where the observed value falls within the distribution of
240 simulated values. This quantile score is equivalent to the p -value of a one-sided statistical test. In previous
241 CSEP experiments, critical values of $\alpha = 0.01$ or $\alpha = 0.05$ were commonly chosen to reject the null

242 hypothesis that the forecast could have generated the observed locations. However, in practice, we use
 243 the consistency tests as diagnostic tools to indicate a degree of (dis)agreement between a forecast and
 244 observations during the testing period (e.g., Bayona et al., 2022). Fig. 4b shows the S-test evaluation for
 245 time-independent Italian forecasts (originally published by Taroni et al., 2018).

246 **Magnitude test** The magnitude (M) test assesses the null hypothesis that the observed magnitude
 247 distribution is consistent with that of the forecast. Similarly to the S-test, the M-test (Zechar et al.,
 248 2010) first sums rates in each magnitude bin over spatial cells and normalizes the forecast so that N_{fore}
 249 matches N_{obs} , thus isolating the magnitude distribution of the forecast. As with the S-test, the M-test then
 250 determines the quantile of the observed (magnitude) joint log-likelihood score in the distribution of joint
 251 log-likelihood scores simulated from the forecast. Observed scores in the tail of the model distribution
 252 indicate discrepancies between the forecast and data that might be scientifically interesting.

253 **Conditional likelihood test** The conditional likelihood (cL) test (Werner et al., 2010, 2011) null
 254 hypothesis states that the observed locations and magnitudes are consistent with the forecast conditional
 255 on the number of observed earthquakes, i.e. the test checks the joint space-magnitude distribution against
 256 the forecast. First, one computes the observed joint log-likelihood score by summing bin-wise log-
 257 likelihood scores over all space-magnitude bins. In this evaluation, the forecast rates are not normalized
 258 to match the observed rate. Again, we assess where this score falls in the critical range of the simulated
 259 distribution of joint log-likelihood scores. Small quantile scores again indicate discrepancies. Effectively,
 260 the CL test represents a combination of the S and M tests.

Comparative testing pyCSEP also provides comparative T- and W-tests (Rhoades et al., 2011) to
 evaluate the relative performance of two models, based on information gain scores per earthquake:

$$IGPE = \frac{1}{N} \sum_{i=1}^N [X_i - Y_i] - \frac{N_A - N_B}{N}, \quad (1)$$

261 where N is the number of observed earthquakes, and $X_i = \ln A(k_i)$ and $Y_i = \ln B(k_i)$ are the log-likelihood
 262 scores obtained by model A and model B in the bin k in which earthquake i occurred, and N_A and N_B
 263 are the expected number of earthquakes according to forecast A and B, respectively. The T-test assesses
 264 whether the $IGPE$ is statistically different from zero. Following Rhoades et al. (2011), one applies the
 265 Student's t-test to the $IGPE$ score of forecast A over forecast B. We consider forecast A to be significantly
 266 more skillful than forecast B if the $IGPE$ is positive and the confidence interval based on the Student's
 267 t-distribution does not include zero. Conversely, if the $IGPE$ is negative and the confidence interval does
 268 not include zero, forecast B is significantly more informative than model A. If the confidence interval
 269 includes zero, we consider differences in the score to be statistically insignificant. Fig. 5 shows T-test

270 results for Californian and Italian forecasts, which were originally published by Zechar et al. (2013) and
271 Taroni et al. (2018), respectively.

272 **Testing catalog-based forecasts** For catalog-based forecasts, pyCSEP provides (1) a number (N) test
273 that compares the (non-Poissonian) number distribution from the forecasts against the observed number
274 of earthquakes; (2) a magnitude (M) test that compares the sum of bin-wise differences in the incremental
275 magnitude distribution; (3) a spatial (S) test that compares the geometric mean of the target event rates; and
276 4) a pseudo-likelihood test based on a statistic that resembles the likelihood of a continuous point-process
277 (Savran et al., 2020). These tests are essentially analogues of the aforementioned consistency tests, but
278 they relax the Poissonian assumption. For a full description of these evaluations and their application to
279 UCERF3-ETAS forecasts made during the 2019 Ridgecrest earthquake sequence in California, see Savran
280 et al. (2020). In Fig. 6, we show an example of the N-test and S-test for a single seven-day UCERF3-ETAS
281 forecast made immediately after the occurrence of the M7.1 mainshock of the Ridgecrest sequence. The
282 catalog-based evaluations are available in the *catalog_evaluations* module in pyCSEP.

283 **Plotting and Other Utilities**

284 Along with the routines for statistical tests, pyCSEP provides a thin wrapper around the matplotlib (Hunter,
285 2007) and cartopy (Met Office, 2015) plotting libraries to provide functions that visualize test results,
286 catalogs, and spatial forecast maps (Fig. 7). We aim to keep the plotting capabilities both easily accessible
287 for early users (i.e., by calling simple methods within most of pyCSEP core classes) and customizable
288 enough to provide journal-quality figures, including: text formatting, legend and colormap editing, spatial
289 grids, and preparing multi-panel figures. The implementation provides access to cartopy’s projection
290 capabilities as well as basic maps, along with various (or user-defined) web-service tiled maps. We intend
291 to keep the plotting functions modular, so that multiple outputs can be combined in single figures, and to
292 preserve the plots if the user requires post-processing of the data or results (as shown in Fig. 7).

293 **REPRODUCIBILITY PACKAGES**

294 CSEP forecasting experiments have run in testing centers, which provide a controlled environment that
295 prevents any access and modification of ongoing experiments. Because pyCSEP now provides the ability
296 to configure bespoke earthquake forecasting experiments, we anticipate that researchers will be interested
297 in using these methods to evaluate their own forecasts. We encourage researchers that use pyCSEP in
298 their publications to follow the approach outlined by Krafczyk et al. (2021) and provide a reproducibility
299 package for their publication.

300 A reproducibility package is a structured set of code, data, and other files that are required to recreate
301 all figures and tables within a manuscript. To illustrate this principle, we provide a reproducibility

302 package for this manuscript. The entry point of the reproducibility package is a script with the following
303 responsibilities: (1) retrieve and verify data artifacts from Zenodo; (2) create a Docker image with the
304 version of pyCSEP, and its dependencies, used for this publication; and (3) run a program to reproduce the
305 figures from this article. Once the reader obtains the reproducibility package, there is a single command
306 to reproduce all of the figures from this paper. We encourage users to try and run the reproducibility
307 package for this manuscript (link for the reproducibility package in the *Data and Resources* section).

308 **PYCSEP COMMUNITY**

309 The pyCSEP efforts aim to strengthen the community of earthquake scientists with an interest in fore-
310 casting. We intend to unite researchers interested in all aspects of earthquake forecasting from model
311 development to testing and evaluation to make the process of forecast testing as transparent and accessible
312 as possible. Fig. 8 is a screenshot from our first community workshop, held (virtually) in March 2021
313 for modelers involved in the project RISE (*‘Real-time earthquake risk reduction for a resilient Europe’*,
314 financed by the European Commission’s *Horizon 2020* program. The workshop introduced forecast
315 developers to the pyCSEP toolkit, helped to identify where improvements and extensions could be made,
316 and invited modelers to contribute. It was held over three sessions, with the first introducing pyCSEP
317 testing, the second allowing modelers to present their current forecasting work, and the third focusing on
318 a hands-on tutorial session. The workshop brought together modelers and model testers to understand the
319 needs of both groups and familiarize all participants with the testing and visualization options currently
320 available in the toolkit. This was later followed by a workshop on contributing to the pyCSEP project
321 through GitHub to familiarize interested users with open-source community software development.

322 Two tutorials were created for the workshop to demonstrate the process of model testing with pyCSEP
323 for grid-based and catalog-based forecasts. The tutorials are in the form of interactive Jupyter notebooks
324 (Kluyver et al., 2016) that provide a template for the key steps of model testing with pyCSEP. Both
325 tutorials use real forecasts and catalog data similar to the examples in this paper. The tutorials are
326 available on the pyCSEP online documentation (link in Data and Resources section), which also includes
327 an installation guide, and a detailed user guide that covers the core concepts to details need to extend
328 pyCSEP functionality.

329 **Open Call for Developers**

330 The workshops highlighted that pyCSEP greatly benefits from an active engagement of its community.
331 Sustaining the development is a community effort and new contributions are essential to extend and
332 improve pyCSEP’s utility. In this regard—and to leverage the open-source development approach—we
333 welcome researchers and developers to join our community and to contribute new ideas and methods

334 (e.g., advanced evaluation capabilities, more robust tests, more efficient testing, etc.). Within the GitHub
335 repository, these contributions can be introduced in the form of ‘pull requests’ (i.e., suggested code
336 changes, improved documentation), or ‘issues’ (e.g., comments or suggestions about technical and
337 scientific approaches). The contributions are transparent and the community can discuss them together.
338 The pyCSEP community additionally meets in regular (developer) calls to coordinate contributions more
339 interactively (e.g., by reviewing source code and new ideas).

340 **CONCLUSION**

341 pyCSEP is an open-source Python package that provides routines for evaluating probabilistic earthquake
342 forecasting models that are expressed as earthquake rates in discrete space-magnitude cells and simulation-
343 based forecasts consisting of synthetic earthquake catalogs. pyCSEP also includes utilities for visualizing
344 forecasts and earthquake catalogs, and configuring earthquake forecasting experiments. The implementa-
345 tion follows best-practices for open-source software development including documentation and continuous
346 integration to build and test new code contributions. In CSEP, we are adopting a software development
347 process that encourages contributions from researchers. To date we have received contributions that have
348 added new evaluation methods and improved plotting capabilities. We advocate that publications involving
349 pyCSEP are accompanied by reproducibility packages. Additionally, we have started a workshop series to
350 train researchers on using pyCSEP and collaborating in open-source development. In 2021, we hosted
351 two workshops teaching users how to use pyCSEP and to work collaboratively in GitHub. We encourage
352 all interested users to visit the online documentation and the code repository to learn more about pyCSEP.

353 **DATA AND RESOURCES**

354 The pyCSEP software can be found on GitHub at <https://github.com/SCECCode/pycsep>
355 and the documentation can be found at <https://docs.cseptest.org>. The reproducibility
356 package for this manuscript can be found at <https://doi.org/10.5281/zenodo.6626265>
357 and the data can be found at <https://doi.org/10.5281/zenodo.5777992>. Best-practices
358 from Computational Infrastructure for Geodynamics (CIG) can be found at <https://geodynamics.org/software/software-bp>. The link to GitHub actions documentation can be found at <https://docs.github.com/en/actions>. The RISE project website can be found at <http://www.rise-eu.org>. Map-tiles for plotting maps can be found at <https://maps.stamen.com>. All
361 websites were last accessed on 24 June 2022.

363 **ACKNOWLEDGMENTS**

364 We would like to thank the editors and two anonymous reviewers for their comments on this article. This
365 research was supported by the Southern California Earthquake Center (Contribution No. 11740). SCEC is
366 funded by NSF Cooperative Agreement EAR-1600087 & USGS Cooperative Agreement G17AC00047.
367 This project has received funding from the European Union's *Horizon 2020 research and innovation*
368 *program* under Grant Agreement Number 821115, *Real-Time Earthquake Risk Reduction for a Resilient*
369 *Europe* (RISE).

370 REFERENCES

- 371 Anzt, H., Bach, F., Druskat, S., Löffler, F., Loewe, A., Renard, B. Y., Seemann, G., Struck, A., Achhammer,
372 E., Aggarwal, P., Appel, F., Bader, M., Brusch, L., Busse, C., Chourdakis, G., Dabrowski, P. W., Ebert,
373 P., Flemisch, B., Friedl, S., Fritsch, B., Funk, M. D., Gast, V., Goth, F., Grad, J.-N., Hegewald, J.,
374 Hermann, S., Hohmann, F., Janosch, S., Kutra, D., Linxweiler, J., Muth, T., Peters-Kottig, W., Rack, F.,
375 Raters, F. H., Rave, S., Reina, G., Reißig, M., Ropinski, T., Schaarschmidt, J., Seibold, H., Thiele, J. P.,
376 Uekermann, B., Unger, S., and Weeber, R. (2021). An environment for sustainable research software in
377 germany and beyond: current state, open challenges, and call for action. *F1000Research*, 9:295.
- 378 Baker, M. (2016). Reproducibility crisis. *Nature*, 78(26):353–66.
- 379 Bayona, J., Savran, W., Strader, A., Hainzl, S., Cotton, F., and Schorlemmer, D. (2021). Two global
380 ensemble seismicity models obtained from the combination of interseismic strain measurements and
381 earthquake-catalogue information. *Geophysical Journal International*, 224(3):1945–1955.
- 382 Bayona, J. A., Savran, W. H., Rhoades, D. A., and Werner, M. J. (2022). Prospective evaluation of
383 multiplicative hybrid earthquake forecasting models in California. *Geophysical Journal International*.
384 ggac018.
- 385 Bird, P., Jackson, D. D., Kagan, Y. Y., Kreemer, C., and Stein, R. S. (2015). Gear1: A global earthquake
386 activity rate model constructed from geodetic strain rates and smoothed seismicity. *Bulletin of the*
387 *Seismological Society of America*, 105(5):2538–2554.
- 388 Bird, P. and Liu, Z. (2007). Seismic hazard inferred from tectonics: California. *Seismological Research*
389 *Letters*, 78(1):37–48.
- 390 Dekel, E. and Feinberg, Y. (2006). Non-bayesian testing of a stochastic prediction. *The Review of*
391 *Economic Studies*, 73(4):893–906.
- 392 Ebel, J. E., Chambers, D. W., Kafka, A. L., and Baglivo, J. A. (2007). Non-poissonian earthquake
393 clustering and the hidden markov model as bases for earthquake forecasting in california. *Seismological*
394 *Research Letters*, 78(1):57–65.
- 395 Eberhard, D. A. J., Zechar, J. D., and Wiemer, S. (2012). A prospective earthquake forecast experiment in
396 the western pacific. *Geophysical Journal International*, 190(3):1579–1592.
- 397 Field, E. H. (2007). Overview of the working group for the development of regional earthquake likelihood
398 models (reim). *Seismological Research Letters*, pages 1–10.
- 399 Geller, R. J. (1997). Earthquake prediction: a critical review. *Geophysical Journal International*,
400 131(3):425–450.
- 401 Gordon, J. S., Clements, R. A., Schoenberg, F. P., and Schorlemmer, D. (2015). Voronoi residuals and
402 other residual analyses applied to csep earthquake forecasts. *Spatial Statistics*, 14:133–150.

403 Group, M. W. (2004). Redazione della mappa di pericolosità sismica prevista dall'ordinanza pc del 20
404 marzo 2003, n. 3274, all. 1 rapporto conclusivo.

405 Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser,
406 E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M. H., Brett, M.,
407 Haldane, A., del Río, J. F., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K., Reddy, T.,
408 Weckesser, W., Abbasi, H., Gohlke, C., and Oliphant, T. E. (2020). Array programming with NumPy.
409 *Nature*, 585(7825):357–362.

410 Helmstetter, A., Kagan, Y. Y., and Jackson, D. D. (2007). High-resolution time-independent grid-based
411 forecast for $m \geq 5$ earthquakes in california. *Seismological Research Letters*, 78(1):78–86.

412 Hunter, J. D. (2007). Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*,
413 9(3):90–95.

414 Jordan, T. H., Chen, Y. T., Gasparini, P., Madariaga, R., Main, I., Marzocchi, W., and Papadopoulos, G.
415 (2011). Operational earthquake forecasting. state of knowledge and guidelines for utilization. *Annals*
416 *of Geophysics*.

417 Jordan, T. H. and Jones, L. M. (2010). Operational earthquake forecasting: Some thoughts on why and
418 how. *Seismological Research Letters*, 81(4):571–574.

419 Kluyver, T., Ragan-Kelley, B., Pérez, F., Granger, B., Bussonnier, M., Frederic, J., Kelley, K., Hamrick, J.,
420 Grout, J., Corlay, S., Ivanov, P., Avila, D., Abdalla, S., and Willing, C. (2016). Jupyter notebooks – a
421 publishing format for reproducible computational workflows. In Loizides, F. and Schmidt, B., editors,
422 *Positioning and Power in Academic Publishing: Players, Agents and Agendas*, pages 87 – 90. IOS
423 Press.

424 Krafczyk, M. S., Shi, A., Bhaskar, A., Marinov, D., and Stodden, V. (2021). Learning from reproducing
425 computational results: introducing three principles and the reproduction package. *Philosophical*
426 *Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 379(2197).

427 Lombardi, A. and Marzocchi, W. (2010a). The assumption of poisson seismic-rate variability in csep/reln
428 experiments. *Bulletin of the Seismological Society of America*, 100(5A):2293–2300.

429 Lombardi, A. M. and Marzocchi, W. (2010b). The etas model for daily forecasting of italian seismicity in
430 the csep experiment. *Annals of Geophysics*.

431 Marzocchi, W., Garcia-Aristizabal, A., Gasparini, P., Mastellone, M. L., and Di Ruocco, A. (2012). Basic
432 principles of multi-risk assessment: a case study in italy. *Natural hazards*, 62(2):551–573.

433 Marzocchi, W., Lombardi, A. M., and Casarotti, E. (2014). The establishment of an operational earthquake
434 forecasting system in italy. *Seismological Research Letters*, 85(5):961–969.

435 McKinney, W. (2010). Data Structures for Statistical Computing in Python. In Stéfan van der Walt and

436 Jarrod Millman, editors, *Proceedings of the 9th Python in Science Conference*, pages 56 – 61.

437 Met Office (2010 - 2015). *Cartopy: a cartographic python library with a Matplotlib interface*. Exeter,
438 Devon.

439 Molchan, G., Romashkova, L., and Peresan, A. (2017). On some methods for assessing earthquake
440 predictions. *Geophysical Journal International*, 210(3):1474–1480.

441 Nandan, S., Ouillon, G., Sornette, D., and Wiemer, S. (2019). Forecasting the full distribution of
442 earthquake numbers is fair, robust, and better. *Seismological Research Letters*, 90(4):1650–1659.

443 National Academies of Sciences, Engineering, and Medicine and others (2019). *Reproducibility and
444 replicability in science*. National Academies Press.

445 pandas Development Team, T. (2020). pandas-dev/pandas: Pandas.

446 pyCSEP Developers (2021). pycsep online documentation.

447 Rhoades, D. A., Christophersen, A., Gerstenberger, M. C., Liukis, M., Silva, F., Marzocchi, W., Werner,
448 M. J., and Jordan, T. H. (2018). Highlights from the first ten years of the new zealand earthquake
449 forecast testing center. *Seismological Research Letters*, 89(4):1229–1237.

450 Rhoades, D. A., Schorlemmer, D., Gerstenberger, M. C., Christophersen, A., Zechar, J. D., and Imoto, M.
451 (2011). Efficient testing of earthquake forecasting models. *Acta Geophysica*, 59(4):728–747.

452 Savran, W. H., J., W. M., D., S., and J., M. P. (2022). pycsep: A python toolkit for earthquake forecast
453 developers. *Journal of Open Source Software*.

454 Savran, W. H., Werner, M. J., Marzocchi, W., Rhoades, D. A., Jackson, D. D., Milner, K., Field, E., and
455 Michael, A. (2020). Pseudoprospective evaluation of ucerf3-etas forecasts during the 2019 ridgecrest
456 sequence. *Bulletin of the Seismological Society of America*, 110(4):1799–1817.

457 Schneider, M., Clements, R., Rhoades, D., and Schorlemmer, D. (2014). Likelihood-and residual-
458 based evaluation of medium-term earthquake forecast models for california. *Geophysical Journal
459 International*, 198(3):1307–1318.

460 Schorlemmer, D., Gerstenberger, M., Wiemer, S., Jackson, D. D., and Rhoades, D. A. (2007). Earthquake
461 likelihood model testing. *Seismological Research Letters*, 78(1):17–29.

462 Schorlemmer, D. and Gerstenberger, M. C. (2007). Relm testing center. *Seismological Research Letters*,
463 78(1):30–36.

464 Schorlemmer, D., Werner, M. J., Marzocchi, W., Jordan, T. H., Ogata, Y., Jackson, D. D., Mak, S.,
465 Rhoades, D. A., Gerstenberger, M. C., Hirata, N., Liukis, M., Maechling, P. J., Strader, A., Taroni, M.,
466 Wiemer, S., Zechar, J. D., and Zhuang, J. C. (2018). The collaboratory for the study of earthquake
467 predictability: Achievements and priorities. *Seismological Research Letters*, 89(4):1305–1313.

468 Stodden, V., Seiler, J., and Ma, Z. (2018). An empirical analysis of journal policy effectiveness for

469 computational reproducibility. *Proceedings of the National Academy of Sciences*, 115(11):2584–2589.

470 Strader, A., Werner, M., Bayona, J., Maechling, P., Silva, F., Liukis, M., and Schorlemmer, D. (2018).

471 Prospective evaluation of global earthquake forecast models: 2 yrs of observations provide preliminary

472 support for merging smoothed seismicity with geodetic strain rates. *Seismological Research Letters*.

473 Taroni, M., Marzocchi, W., Schorlemmer, D., Werner, M. J., Wiemer, S., Zechar, J. D., Heiniger, L., and

474 Euchner, F. (2018). Prospective csep evaluation of 1-day, 3-month, and 5-yr earthquake forecasts for

475 italy. *Seismological Research Letters*, 89(4):1251–1261.

476 Tsuruoka, H., Hirata, N., Schorlemmer, D., Euchner, F., Nanjo, K. Z., and Jordan, T. H. (2012). Csep

477 testing center and the first results of the earthquake forecast testing experiment in japan. *Earth, Planets*

478 *and Space*, 64(8):661–671.

479 Werner, M. J., Helmstetter, A., Jackson, D. D., and Kagan, Y. Y. (2011). High-resolution long-term

480 and short-term earthquake forecasts for california. *Bulletin of the Seismological Society of America*,

481 101(4):1630–1648.

482 Werner, M. J., Helmstetter, A., Jackson, D. D., Kagan, Y. Y., and Wiemer, S. (2010). Adaptively smoothed

483 seismicity earthquake forecasts for italy. *arXiv preprint arXiv:1003.4374*.

484 Werner, M. J. and Sornette, D. (2008). Magnitude uncertainties impact seismic rate estimates, forecasts,

485 and predictability experiments. *Journal of Geophysical Research: Solid Earth*, 113(B8).

486 Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G., Axton, M., Baak, A., Blomberg, N.,

487 Boiten, J.-W., da Silva Santos, L. B., Bourne, P. E., Bouwman, J., Brookes, A. J., Clark, T., Crosas,

488 M., Dillo, I., Dumon, O., Edmunds, S., Evelo, C. T., Finkers, R., Gonzalez-Beltran, A., Gray, A. J. G.,

489 Groth, P., Goble, C., Grethe, J. S., Heringa, J., 't Hoen, P. A. C., Hooft, R., Kuhn, T., Kok, R., Kok, J.,

490 Lusher, S. J., Martone, M. E., Mons, A., Packer, A. L., Persson, B., Rocca-Serra, P., Roos, M., van

491 Schaik, R., Sansone, S.-A., Schultes, E., Sengstag, T., Slater, T., Strawn, G., Swertz, M. A., Thompson,

492 M., van der Lei, J., van Mulligen, E., Velterop, J., Waagmeester, A., Wittenburg, P., Wolstencroft,

493 K., Zhao, J., and Mons, B. (2016). The FAIR Guiding Principles for scientific data management and

494 stewardship. *Scientific Data*, 3(1):160018.

495 Zechar, J. D., Gerstenberger, M. C., and Rhoades, D. A. (2010). Likelihood-based tests for evaluat-

496 ing space-rate-magnitude earthquake forecasts. *Bulletin of the Seismological Society of America*,

497 100(3):1184–1195.

498 Zechar, J. D. and Jordan, T. H. (2010). Simple smoothed seismicity earthquake forecasts for italy. *Annals*

499 *of Geophysics*, 53(3):99–105.

500 Zechar, J. D., Schorlemmer, D., Werner, M. J., Gerstenberger, M. C., Rhoades, D. A., and Jordan, T. H.

501 (2013). Regional earthquake likelihood models i: First-order results. *Bulletin of the Seismological*

502 *Society of America*, 103(2A):787–798.

503 Zechar, J. D. and Zhuang, J. (2010). Risk and return: evaluating reverse tracing of precursors earthquake
504 predictions. *Geophysical Journal International*, 182(3):1319–1326.

505 **AUTHOR LIST**

506 **William Savran**, Zumberge Hall of Science (ZHS), 3651 Trousdale Pkwy, Los Angeles, CA 90089-0740

507 **José A. Bayona**, School of Earth Sciences, University of Bristol, Wills Memorial Building, Queens Road,
508 Bristol BS8 1RJ, United Kingdom

509 **Pablo Iturrieta**, Helmholtzstraße 7, 14467 Potsdam, Germany

510 **Khawaja M. Asim**, Helmholtzstraße 7, 14467 Potsdam, Germany

511 **Han Bao**, Department of Earth, Planetary, and Space Sciences, University of California, Los Angeles,
512 595 Charles Young Drive East, Box 951567, Los Angeles, CA 90095-1567

513 **Kirsty Bayliss**, The King's Buildings, James Hutton Road, Edinburgh EH9 3FE

514 **Marcuss Hermann**, Corso Umberto I, 40, 80138 Napoli NA, Italy

515 **Danijel Schorlemmer**, Helmholtzstraße 7, 14467 Potsdam, Germany

516 **Philip J. Maechling**, Zumberge Hall of Science (ZHS), 3651 Trousdale Pkwy, Los Angeles, CA 90089-
517 0740

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

520 List of Figures

521 Figure 1: Schematic of the pyCSEP classes and code structure showing the core classes
522 of pyCSEP. For a complete description of the software, please see the online
523 documentation (pyCSEP Developers, 2021).

524 Figure 2: (a) Time-independent grid-based forecasts for the HELMSTETTER model for
525 California (Helmstetter et al., 2007), and (b) the MELETTI model for Italy
526 (Group, 2004). The colormap reflects the logarithm of the expected rate of
527 **M**4.5+ earthquakes in $0.1^\circ \times 0.1^\circ$ spatial bins over a five year period. The red
528 circles depict locations of observed earthquakes during the five-year evaluation
529 period. Earthquakes are shown atop the HELMSTETTER forecast from 01
530 January 2006 through 01 January 2011, and atop the MELETTI forecast from 01
531 January 2010 through 01 January 2015.

532 Figure 3: Select realizations (synthetic catalogs) from a week-long UCERF3-ETAS fore-
533 cast generated during the 2019 Ridgecrest, California, sequence. The forecast
534 starts immediately following the **M**7.1 mainshock. Catalogs are chosen based on
535 their percentile in the forecasted number distribution: (a) shows the 5th percentile,
536 (b) shows the median, (c) shows the 75th percentile, and (d) shows the 99.9th
537 percentile catalog. Individual earthquakes are represented by red circles, and the
538 background image shows the expected rate of **M**2.5+ earthquakes aggregated
539 in $0.025^\circ \times 0.025^\circ$ spatial bins (i.e. the ensemble average over the simulated
540 catalogs).

541 Figure 4: (a) N-test results for the HELMSTETTER (Helmstetter et al., 2007), BIRD.LIU
542 (Bird and Liu, 2007), and EBEL (Ebel et al., 2007) earthquake forecasts for
543 California. The markers depict the number of M4.95+ earthquakes during
544 the 2006–2010 RELM evaluation period. The green square indicates that this
545 number falls within the 95% range of the forecast number distribution (solid
546 bar), whereas red circles indicate inconsistencies between the forecast and ob-
547 servations. Thus, the observed number of earthquakes is consistent with the
548 HELMSTETTER model, whereas the BIRD.LIU and EBEL models overesti-
549 mate seismicity in the region. (b) Results of the S-test for the LOMBARDI
550 (Lombardi and Marzocchi, 2010b), MELETTI (Group, 2004) and WERNER-
551 MI (Werner et al., 2010) forecasts for Italy. The markers represent the spatial
552 joint log-likelihood of each model. The green square indicates that the spatial
553 distribution forecasted by the LOMBARDI model is consistent with the spatial
554 distribution of observed seismicity at a 0.01 significance level. Red circles
555 indicate that the observed locations are inconsistent with the spatial forecasts by
556 the MELETTI and WERNER models. In both panels, significance levels of the
557 test are chosen from the original publication of these results (i.e., Taroni et al.,
558 2018; Zechar et al., 2013).

559 Figure 5: (a) T-test results comparing the BIRD.LIU and EBEL forecasts with the bench-
560 mark HELMSTETTER forecast (horizontal dashed line) in California. Red
561 circles indicate information gains of the forecasts (here both negative) and
562 vertical red bars show 95% confidence intervals. These results indicate that
563 HELMSTETTER is significantly more informative than BIRD.LIU and EBEL
564 during the evaluation period. (b) Results of the T-test for the LOMBARDI and
565 WERNER forecasts with respect to the MELETTI reference forecast (horizontal
566 dashed line), in Italy. White circles display information gains of the forecasts
567 and vertical gray bars represent 95% confidence intervals. The horizontal dashed
568 line falls within both confidence intervals, indicating that these models are as
569 statistically informative as the MELETTI benchmark model.

570 Figure 6: (a) Expected rate of $M_{2.5+}$ events within $0.025^\circ \times 0.025^\circ$ spatial cells from
571 a week-long UCERF3-ETAS forecast. The expected rates are computed by
572 averaging over the simulated catalogs. The red circles show the locations of
573 the 827 events observed during the forecast period. (b) N-test result for this
574 UCERF3-ETAS forecast. The histogram shows the forecast number distribution
575 with the two-sided 5% critical region highlighted red. The number of observed
576 events, ω , is depicted as the dashed line, which falls near the mode of the forecast
577 distribution, indicating consistency. δ_1 equals the fraction of catalogs in the
578 forecast that predict at least as many earthquakes as the observations, and δ_2
579 denotes the fraction of catalogs that contain at most the number of observed
580 earthquakes. (c) S-test for the UCERF3-ETAS forecast. The histogram shows
581 the distribution of simulated scores, computed by assigning each earthquake an
582 expected rate based on the bin-wise values in (a) and computing the geometric
583 mean over each catalog in the forecast. The two-sided 5% critical region is
584 highlighted red. We compute the same statistic from the observed catalog, and
585 show the value, ω , using the dashed line. The score, γ , shows the quantile where
586 the observed value falls in the forecast distribution.

587 Figure 7: Example of the spatial plotting capabilities of pyCSEP through a thin-wrapper
588 over cartopy. (a) A quick basemap can be obtained from the default plotting
589 arguments that uses map tiles by Stamen Design (link in Data and Resources
590 section). (b) On top of the basemap, two post-processed forecasts (the ratio
591 between them over a given range of magnitudes). (c) The observed catalog
592 within the same magnitude range, with auto-scaled symbols according to their
593 magnitudes. These functions are intended to be used with pyCSEP classes and
594 provide a simple way of visualizing spatial earthquake forecasts and catalogs.

595 Figure 8: A screenshot of the participants at the virtual pyCSEP training workshop in
596 March 2021 hosted by the RISE project.

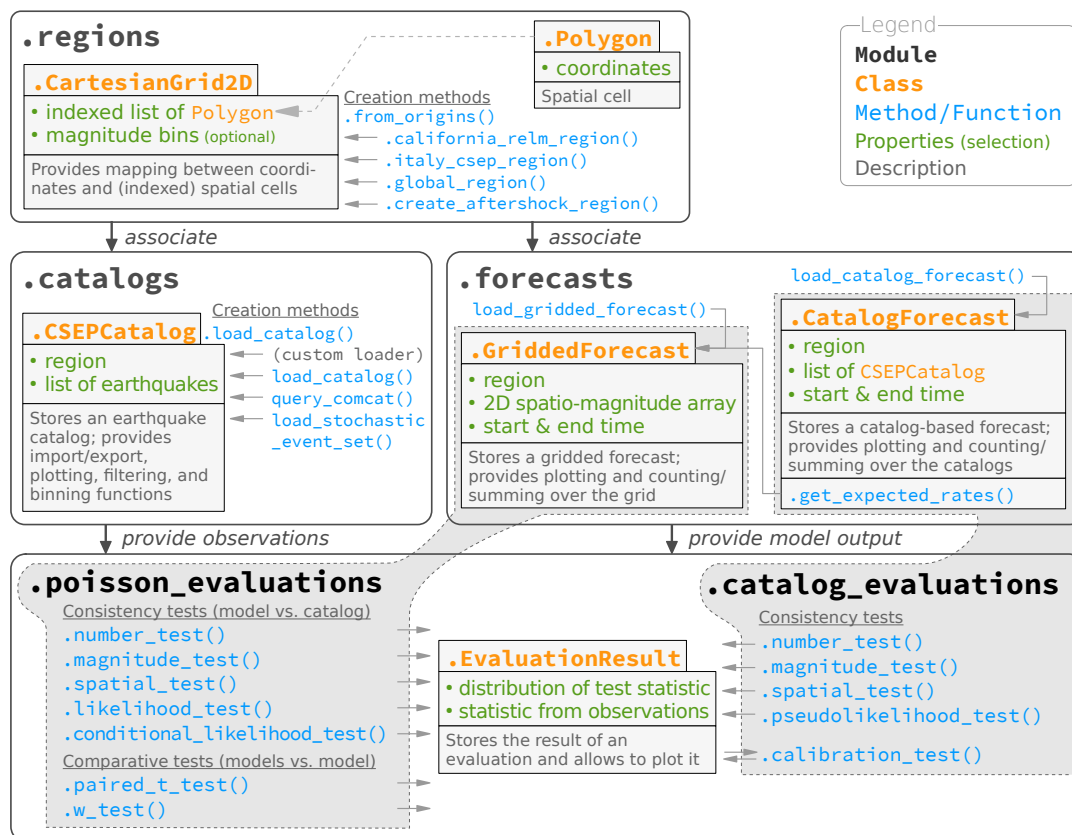


Figure 1. Schematic of the pyCSEP classes and code structure showing the core classes of pyCSEP. For a complete description of the software, please see the online documentation (pyCSEP Developers, 2021).

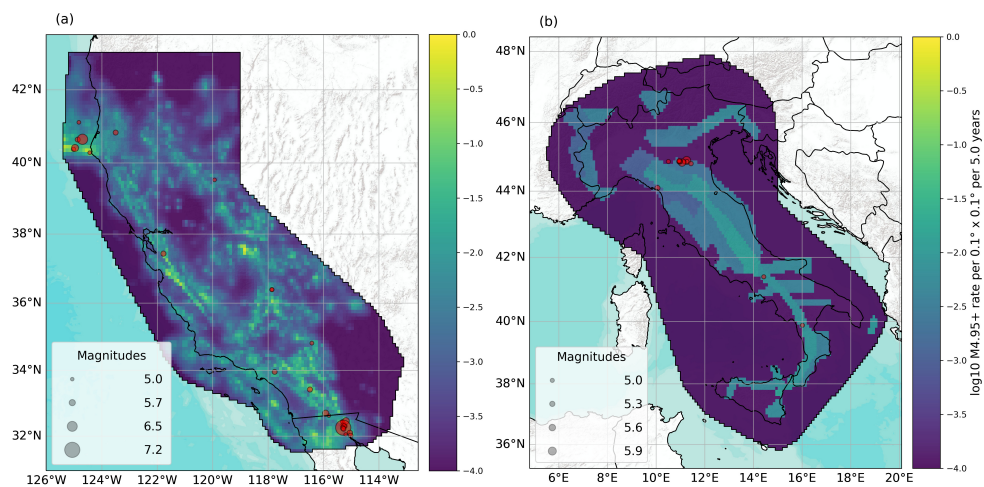


Figure 2. (a) Time-independent grid-based forecasts for the HELMSTETTER model for California (Helmstetter et al., 2007), and (b) the MELETTI model for Italy (Group, 2004). The colormap reflects the logarithm of the expected rate of M_{4.5+} earthquakes in $0.1^\circ \times 0.1^\circ$ spatial bins over a five year period. The red circles depict locations of observed earthquakes during the five-year evaluation period. Earthquakes are shown atop the HELMSTETTER forecast from 01 January 2006 through 01 January 2011, and atop the MELETTI forecast from 01 January 2010 through 01 January 2015.

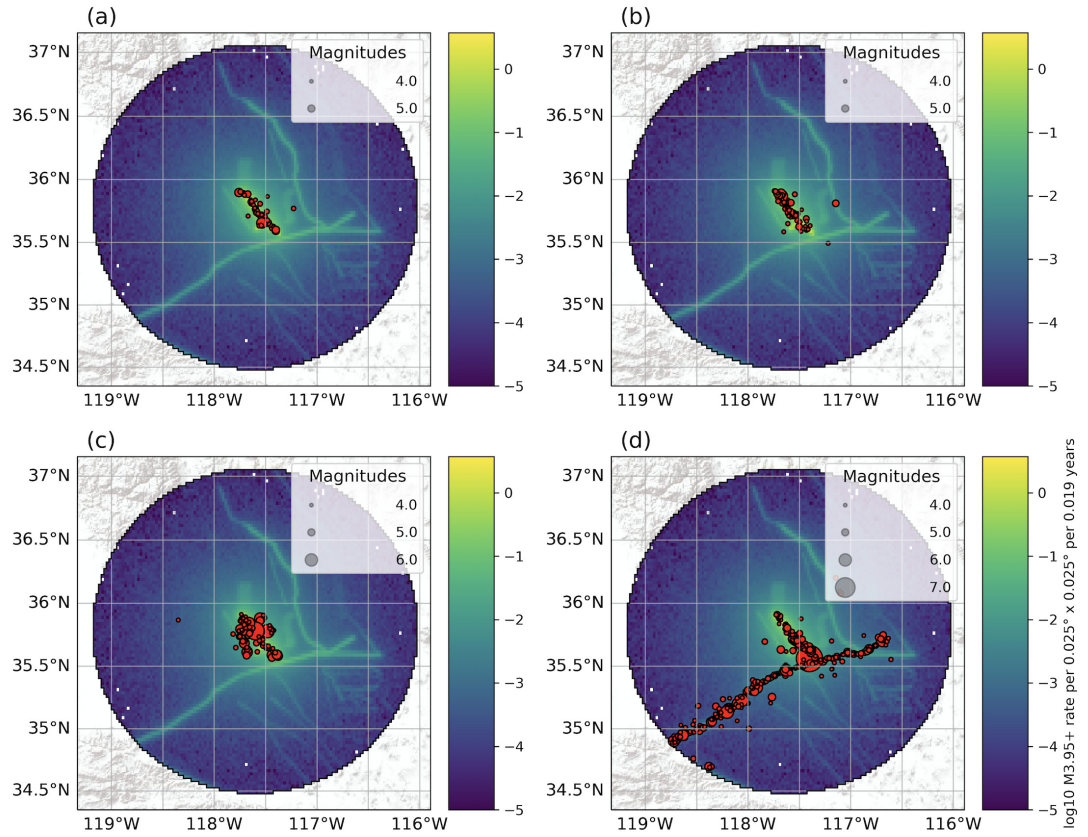


Figure 3. Select realizations (synthetic catalogs) from a week-long UCERF3-ETAS forecast generated during the 2019 Ridgecrest, California, sequence. The forecast starts immediately following the $M_{7.1}$ mainshock. Catalogs are chosen based on their percentile in the forecasted number distribution: (a) shows the 5th percentile, (b) shows the median, (c) shows the 75th percentile, and (d) shows the 99.9th percentile catalog. Individual earthquakes are represented by red circles, and the background image shows the expected rate of $M_{2.5+}$ earthquakes aggregated in $0.025^\circ \times 0.025^\circ$ spatial bins (i.e. the ensemble average over the simulated catalogs).

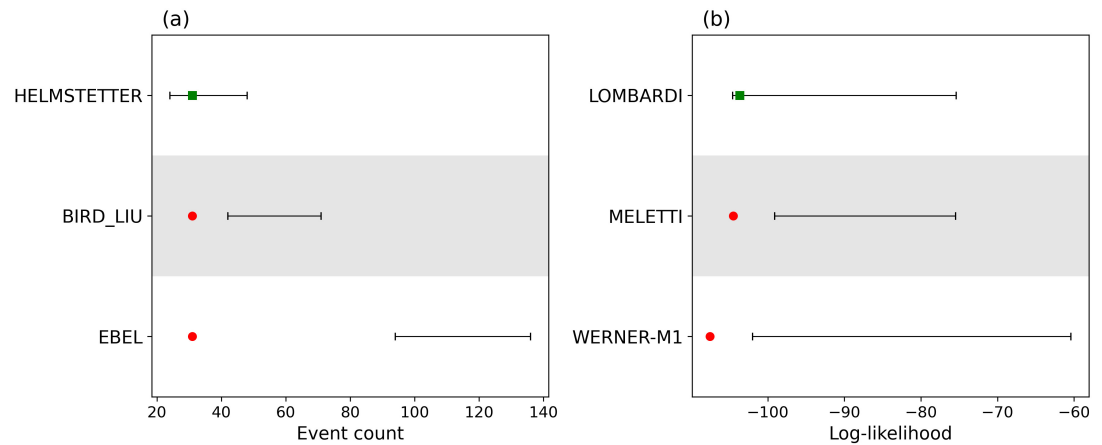


Figure 4. (a) N-test results for the HELMSTETTER (Helmstetter et al., 2007), BIRD_LIU (Bird and Liu, 2007), and EBEL (Ebel et al., 2007) earthquake forecasts for California. The markers depict the number of $M_{4.95+}$ earthquakes during the 2006–2010 RELM evaluation period. The green square indicates that this number falls within the 95% range of the forecast number distribution (solid bar), whereas red circles indicate inconsistencies between the forecast and observations. Thus, the observed number of earthquakes is consistent with the HELMSTETTER model, whereas the BIRD_LIU and EBEL models overestimate seismicity in the region. (b) Results of the S-test for the LOMBARDI (Lombardi and Marzocchi, 2010b), MELETTI (Group, 2004) and WERNER-M1 (Werner et al., 2010) forecasts for Italy. The markers represent the spatial joint log-likelihood of each model. The green square indicates that the spatial distribution forecasted by the LOMBARDI model is consistent with the spatial distribution of observed seismicity at a 0.01 significance level. Red circles indicate that the observed locations are inconsistent with the spatial forecasts by the MELETTI and WERNER models. In both panels, significance levels of the test are chosen from the original publication of these results (i.e., Taroni et al., 2018; Zechar et al., 2013).

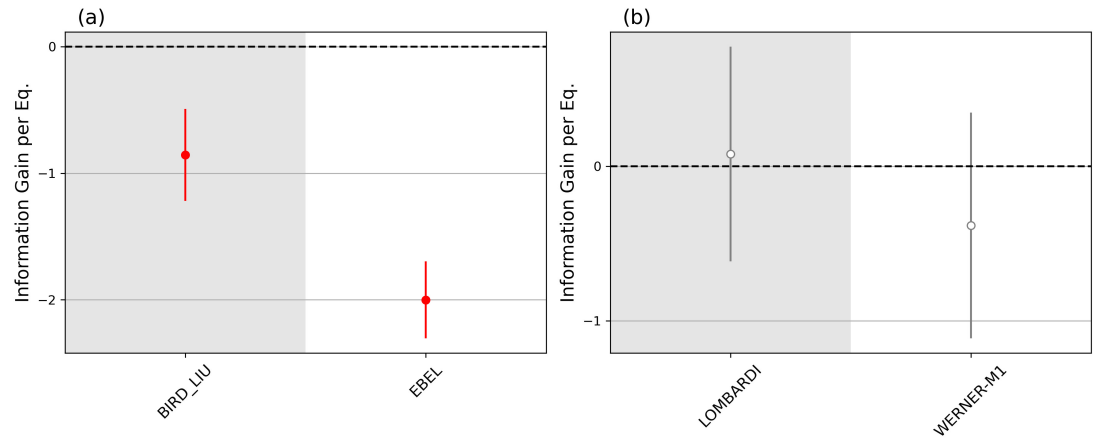


Figure 5. (a) T-test results comparing the BIRD_LIU and EBEL forecasts with the benchmark HELMSTETTER forecast (horizontal dashed line) in California. Red circles indicate information gains of the forecasts (here both negative) and vertical red bars show 95% confidence intervals. These results indicate that HELMSTETTER is significantly more informative than BIRD_LIU and EBEL during the evaluation period. (b) Results of the T-test for the LOMBARDI and WERNER forecasts with respect to the MELETTI reference forecast (horizontal dashed line), in Italy. White circles display information gains of the forecasts and vertical gray bars represent 95% confidence intervals. The horizontal dashed line falls within both confidence intervals, indicating that these models are as statistically informative as the MELETTI benchmark model.

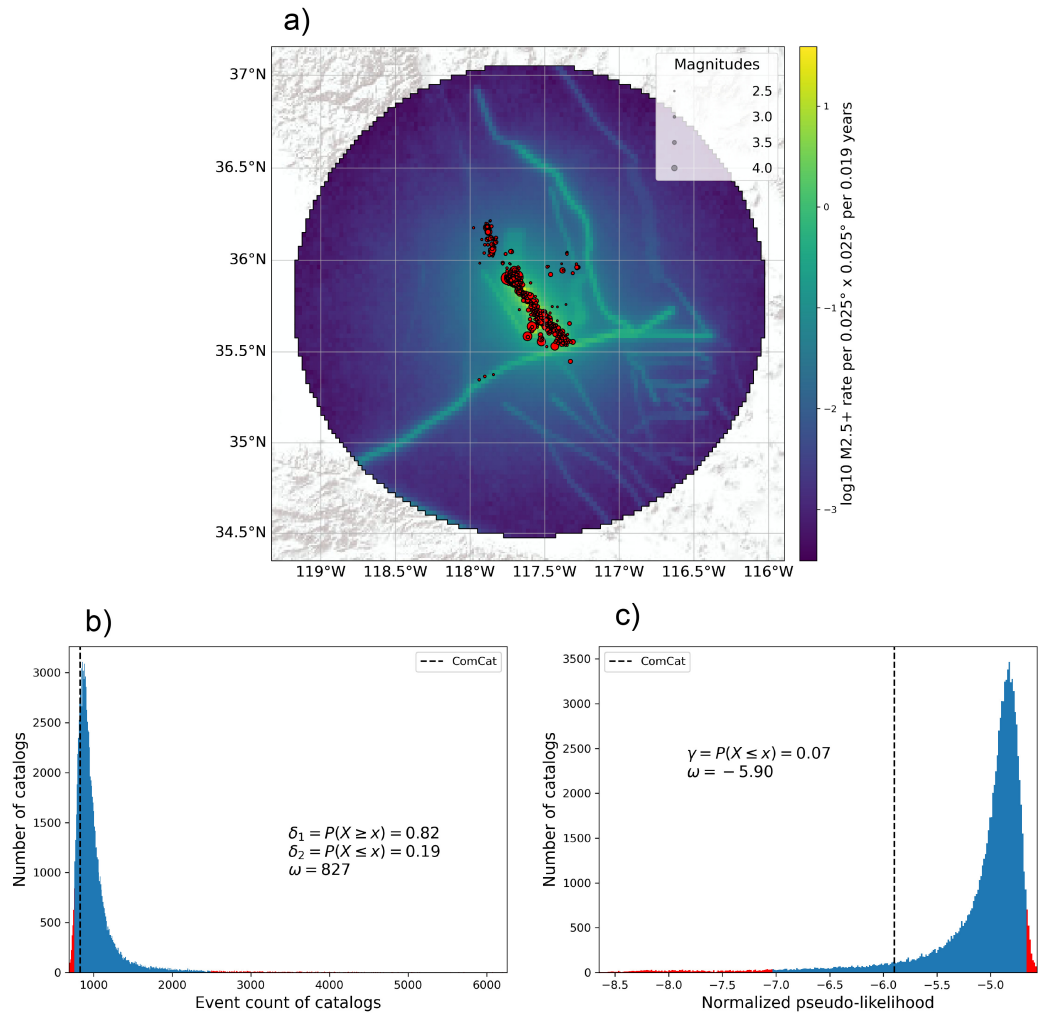


Figure 6. (a) Expected rate of M2.5+ events within $0.025^\circ \times 0.025^\circ$ spatial cells from a week-long UCERF3-ETAS forecast. The expected rates are computed by averaging over the simulated catalogs. The red circles show the locations of the 827 events observed during the forecast period. (b) N-test result for this UCERF3-ETAS forecast. The histogram shows the forecast number distribution with the two-sided 5% critical region highlighted red. The number of observed events, ω , is depicted as the dashed line, which falls near the mode of the forecast distribution, indicating consistency. δ_1 equals the fraction of catalogs in the forecast that predict at least as many earthquakes as the observations, and δ_2 denotes the fraction of catalogs that contain at most the number of observed earthquakes. (c) S-test for the UCERF3-ETAS forecast. The histogram shows the distribution of simulated scores, computed by assigning each earthquake an expected rate based on the bin-wise values in (a) and computing the geometric mean over each catalog in the forecast. The two-sided 5% critical region is highlighted red. We compute the same statistic from the observed catalog, and show the value, ω , using the dashed line. The score, γ , shows the quantile where the observed value falls in the forecast distribution.

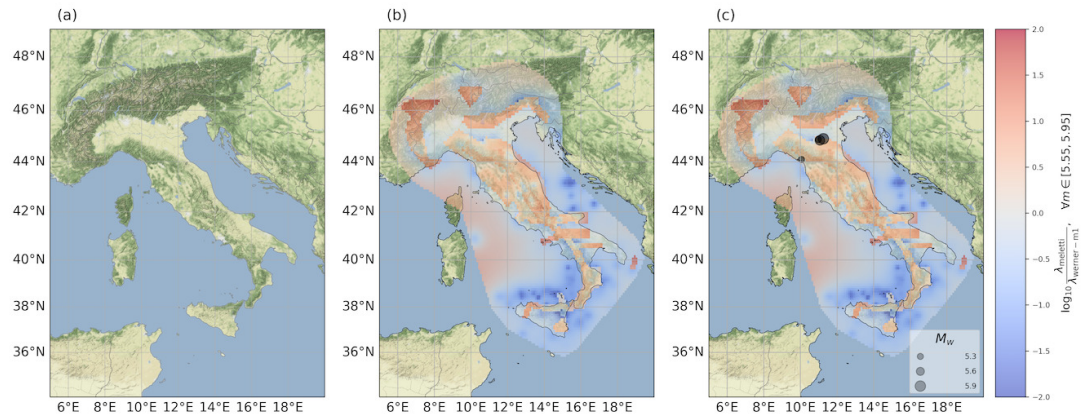


Figure 7. Example of the spatial plotting capabilities of pyCSEP through a thin-wrapper over cartopy. (a) A quick basemap can be obtained from the default plotting arguments that uses map tiles by Stamen Design (link in Data and Resources section). (b) On top of the basemap, two post-processed forecasts (the ratio between them over a given range of magnitudes). (c) The observed catalog within the same magnitude range, with auto-scaled symbols according to their magnitudes. These functions are intended to be used with pyCSEP classes and provide a simple way of visualizing spatial earthquake forecasts and catalogs.

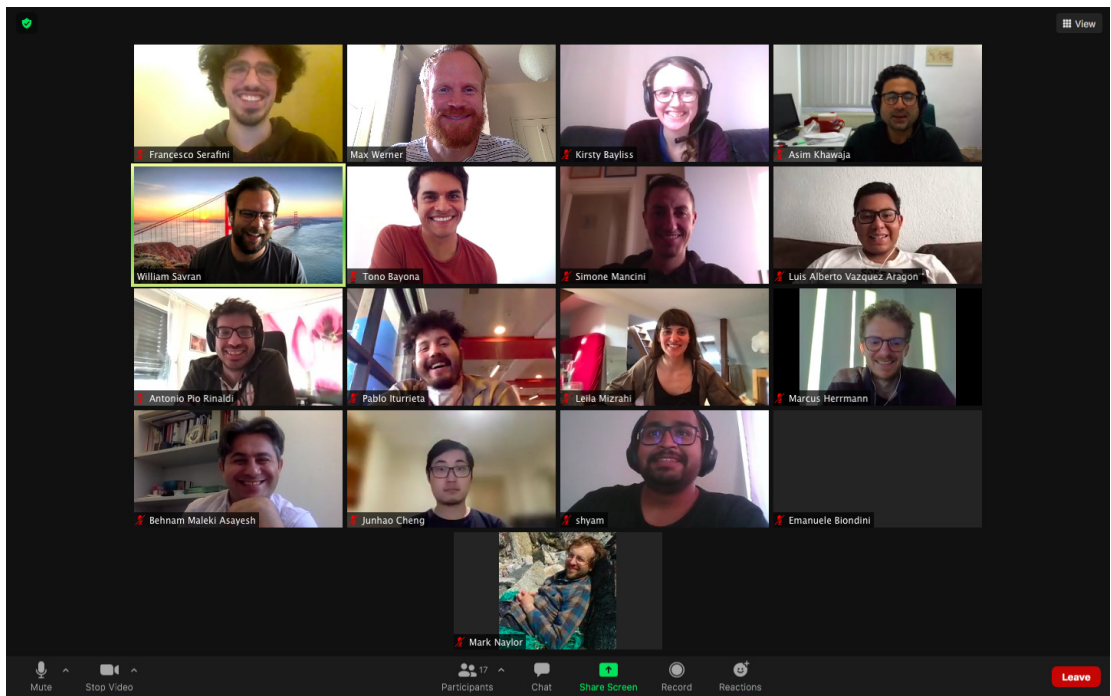


Figure 8. A screenshot of the participants at the virtual pyCSEP training workshop in March 2021 hosted by the RISE project.