

1 **Uncovering Drivers of Atmospheric River Flood Damage using Interpretable** 2 **Machine Learning**

3 Corinne Bowers, Ph.D.^{1,2}, Katherine A. Serafin, Ph.D.³, and Jack W. Baker, Ph.D.¹

4 ¹Civil and Environmental Engineering Department, Stanford University, Stanford, CA.

5 ² [Currently employed at U.S. Geological Survey, Reston, VA. Email: cbowers@usgs.gov.](mailto:cbowers@usgs.gov)

6 ³Department of Geography, University of Florida, Gainesville, FL.

7 **ABSTRACT**

8 The intensity of an atmospheric river (AR) is only one of the factors influencing the damage it
9 will cause. We use random forest models fit to hazard, exposure, and vulnerability data at different
10 spatial and temporal scales in California to predict the probability that a given AR event will
11 cause flood damage, as measured by National Flood Insurance Program (NFIP) claims. We first
12 demonstrate the usefulness of data-driven models and interpretable machine learning to identify
13 and describe drivers of AR flood damage. Hazard features, particularly measures of AR intensity
14 such as total precipitation, increase the probability of damage with increasing values up to a
15 threshold point, after which the probability of damage saturates. While hazard is generally the most
16 important risk dimension across all models, exposure and vulnerability contribute up to a third of the
17 explanatory power. Exposure and variability features generally increase the probability of damage
18 with increasing values, apart from a few instances which can be explained by physical intuition,
19 but tend to affect the probability of damage less for the largest AR events. Comparisons between
20 random forest models at different spatial and temporal scales show general agreement. We then
21 examine limitations inherent in publicly available exposure, vulnerability, and loss data, focusing
22 on the difference in temporal resolution between variables from different risk dimensions and
23 discrepancies between NFIP claims and total flood losses, and describe how those limitations may

24 affect the model results. Overall, the application of interpretable machine learning to understand
25 the contributions of exposure and vulnerability to AR-driven flood risk has identified potential
26 community risk drivers and strategies for resilience, but the results must be considered in the
27 context of the data that produced them.

28 **INTRODUCTION**

29 Flooding is the most common and costly natural disaster that Americans face. Climate change
30 has already increased the frequency and severity of floods; of the 41 billion-dollar flood disasters
31 in the United States since 1980, 18 have occurred in the past decade (NOAA NCEI 2023). Floods
32 become disasters based on not only the intensity of the hazard, but also interactions with the
33 landscape, infrastructure, and communities at a particular location. In order to characterize flood
34 risk, we rely on the well-established definition that risk is the product of three dimensions: hazard,
35 exposure, and vulnerability. Hazard includes the intensity of the atmospheric event as well as
36 environmental factors that could mediate or amplify flooding, such as impervious land cover or wet
37 antecedent conditions. Exposure represents the people and buildings who experience the hazard,
38 and vulnerability quantifies the ability of those people and buildings to withstand the hazard. All
39 three risk dimensions must be accounted for in order to build models of flood damage that are
40 both accurate (able to predict the magnitude of damage expected from a given storm event) and
41 interpretable (able to determine which risk factors contributed most to damage during that event).

42 Flood damage can be modeled using either a process-based or a data-driven approach. Process-
43 based models start with the driving hydroclimatic conditions and simulate the physical processes
44 from streamflow and inundation to damage and loss (e.g., FEMA (2006)). Most process-based
45 models in the literature stop at loss prediction, though, and do not extend their analysis to quantify
46 the drivers of loss. Data-driven models approach the problem from the opposite direction, starting
47 with instances when impacts were observed and working backwards to empirically estimate the
48 factors that are most predictive of impact (Solomatine and Ostfeld 2008). A data-driven model can
49 take many forms, ranging from ordinary least squares regression to complex machine learning (ML)
50 and artificial intelligence (AI) models. ML models have increased in popularity for assessment

51 of flood hazard (Sadler et al. 2018; Mobley et al. 2021) and flood damage (Wagenaar et al. 2017;
52 Szczyrba et al. 2021) because of their high predictive accuracy and the increased availability of
53 data for fitting the models.

54 However, there are two main limitations in existing ML models for flood risk assessment. First,
55 the increased predictive accuracy of more complex model forms comes at a price of decreased model
56 interpretability. Interpretability is especially important in a flood risk context where knowing *why*
57 the model produced a certain outcome matters as much as or more than *what* the outcome was.
58 This limitation has been partially addressed through a suite of tools that fall under the umbrella of
59 “interpretable ML” (Molnar 2023), all developed with the goal of improving the degree to which
60 a human can understand the cause of a data-driven model decision. While some researchers have
61 incorporated interpretable ML results in their flood risk assessments (e.g., Stein et al. (2021)), the
62 practice is not widespread. Second, previous research has shown that data-driven models for flood
63 risk assessment are sensitive to the spatial resolution of the data (Komolafe et al. 2018; Pollack
64 et al. 2022) and that differences in temporal resolution across predictor variables can skew results
65 (Mobley et al. 2021). Very few studies have examined the effect of spatial or temporal resolution on
66 model interpretability results, and to our knowledge none have considered both factors in tandem.

67 In this paper, we build random forest (RF) classification models to predict the likelihood of
68 flood damage due to atmospheric rivers (ARs) in California at different spatial and temporal scales.
69 ARs are the primary drivers of flood risk in the western US, associated with extreme precipitation
70 (Lamjiri et al. 2017), hydrologic floods (Konrad and Dettinger 2017), and economic impacts
71 (Corringham et al. 2019). We create an extensive dataset with over forty predictor variables
72 representing hazard, exposure, and vulnerability. We then use interpretable ML to explore the
73 contributions of these variables to the prevalence of insurance claims from the National Flood
74 Insurance Program (NFIP). Our RF models quantify the value of including information about
75 community-level exposure and social and infrastructural vulnerability in models of AR-driven
76 flood damage and identify nonlinear threshold points and variable interactions that can guide
77 potential resilience strategies. This paper also makes a more general methodological contribution

78 to the literature on ML in flood risk by comparing predictive accuracy and model interpretability
79 results from RF models created at multiple spatial and temporal scales. We offer a perspective on
80 the benefits and limitations of existing publicly available exposure, vulnerability, and loss data, and
81 conclude by proposing avenues of work to improve future data-driven models of both flood risk
82 and flood risk drivers.

83 **DATA**

84 **Response Variable**

85 The response variable of interest is a binary indicator of whether or not an AR storm caused
86 flood damage in a specific geographic unit (county or census tract). We define a damaging
87 AR event as one that causes flood insurance claims to be submitted by a policyholder in the
88 Federal Emergency Management Agency (FEMA) National Flood Insurance Program (NFIP)
89 (FEMA 2023b). We include denied claims and claims below the policy deductible (zero payout)
90 in our analysis, assuming a filed claim indicates that the policyholder experienced some negative
91 consequence due to an AR event. We do not distinguish between pluvial, fluvial, coastal, or
92 indirect flood effects. The number of claims needed to qualify an AR event as damaging depends
93 on the spatial resolution. At the census tract level, the threshold is one claim. At the county
94 level, large differences in population between counties mean that more populous counties have
95 more policyholders and are consequently more likely to have at least one claim filed somewhere
96 in the county during the AR event; we therefore define the threshold as $(N_{county})_i / N_{state}$, where
97 $(N_{county})_i$ is the number of NFIP policies in county i and N_{state} is the statewide median of policies
98 per county. This policy-based threshold corrects the population bias at the county level and more
99 evenly distributes damaging events across the state.

100 NFIP claims are often used as a proxy for flood impacts because claims are available at the
101 census tract scale and tagged to a specific date of loss, which allows for a granular examination of
102 flood impacts. Multiple other researchers have used NFIP claims to fit data-driven models of flood
103 hazard (Mobley et al. 2021) and loss (Czajkowski et al. 2017; Knighton et al. 2020). However, NFIP
104 policyholders are not a representative sample of California residents. Only about 2% of eligible

105 homeowners are insured (FEMA 2023b), and due to a combination of both self-selection within
106 risky areas and the mandatory purchase requirement for homes with federally backed mortgages
107 within the 100-year floodplain, NFIP policyholders are more likely than the general population to
108 live in high-risk areas with a history of flooding (Bradt et al. 2021). Insurance take-up rates are also
109 influenced by education (Atreya et al. 2015), community flood protection investment (Zahran et al.
110 2009), and income and home value (Darlington and Yiannakoulis 2022), among other factors.

111 Previous works have addressed the representation bias of the NFIP in a number of ways, from
112 applying correction factors (Smith and Katz 2013; Corringham and Cayan 2019) to modeling
113 damage at uninsured properties (Thomson et al. 2023). We address it here by limiting our analysis
114 to classification rather than regression. Focusing on damage versus no-damage and neglecting
115 claim payout values avoids issues arising from differences in coverage limits between policyholders
116 within and outside of the 100-year floodplain, coverage limit changes over time, and concerns about
117 overrepresentation of higher-valued properties. However, it means our results will only show the
118 underlying drivers of damage probability, which may or may not be the same as drivers of damage
119 magnitude (Rözer et al. 2019). There also still remain demographic and socioeconomic differences
120 at the intra-county level between who is insured and who is not.

121 **Predictor Variables**

122 Table 1 lists all of the predictors in the model by risk dimension and by concept, where concepts
123 represent groups of related variables. Table 1 also includes references to the data source for each
124 variable as well as references that support that variable's potential connection to the response.
125 Spatial variation of the data is at the census tract scale or smaller (T), at the county scale (C), or
126 constant (-). Temporal variation is at the event level (E), monthly (M), yearly (Y), or constant (-).

127 We provide additional context around the variables chosen to represent each risk dimension,
128 starting with hazard. Each record in our dataset represents one AR event. To identify ARs, we
129 used the Rutz et al. (2014) algorithm, which defines ARs as contiguous areas greater than 2,000
130 km in length and with integrated water transport (IVT) values over 250 kg/m/s. IVT was calculated

131 from the Modern-Era Retrospective analysis for Research and Applications, Version 2 (MERRA-2;
132 [Gelaro et al. \(2017\)](#)). MERRA-2 reports data at a resolution of $0.5^\circ \times 0.625^\circ$ ($\sim 50\text{km} \times 50\text{km}$)
133 from 1980 to present. We recorded the maximum IVT and duration of each AR event and used
134 the [Ralph et al. \(2019\)](#) scale to categorize the intensity of each AR from 1 (mostly beneficial) to 5
135 (mostly hazardous). Antecedent conditions are measured with three variables: total precipitation in
136 the 3 days prior to the AR event, total precipitation in the 14 days prior to the AR event, and average
137 soil moisture over the 3 days prior to the AR event. Large-scale climate modes such as ENSO and
138 PDO capture time periods when flood risk increases or decreases over a broad geographic range,
139 and land surface variables capture on-the-ground conditions that can amplify storm effects.

140 For exposure, we focus on variables related to population and housing. The large majority
141 of NFIP policies cover residential buildings, so NFIP claims are more representative of housing
142 exposure than other types of infrastructure. We do not measure other types of exposed assets such
143 as roads, critical infrastructures, crops and livestock, and cultural heritage sites. We also include
144 variables identifying specific geographies associated with higher NFIP insurance take-up rates.

145 For social vulnerability, we rely on constructed indices, particularly the Centers for Disease
146 Control (CDC) Social Vulnerability Index (SVI) ([Flanagan et al. 2011](#)) and CalEnviroScreen 4.0, a
147 statewide screening tool for vulnerability to environmental hazards ([August et al. 2021](#)). The CDC
148 SVI was chosen because of its long time record, with values extending back to 2000, and its ability
149 to explain both recorded damages and fatalities in an empirical validation exercise ([Bakkensen et al.](#)
150 [2017](#)). We included the four components of SVI (socioeconomic status, household characteristics,
151 racial & ethnic minority status, and housing type & transportation) as separate predictor variables.
152 The CalEnviroScreen metrics (population characteristics, pollution burden, and disadvantaged
153 communities) were chosen because of their calibration to California, their relevance in statewide
154 planning decisions, and their focus on environmental justice. One drawback of constructed indices,
155 though, is that they are designed for comparison over space rather than over time. While the indices
156 are updated regularly, each data generation is normalized such that the values at any given time
157 represent only the relative ranking of one census tract or county, so the values do not necessarily

158 capture absolute changes in vulnerability over time (Bakkensen et al. 2017). Therefore we include
159 median household income, percent of the population as non-Hispanic white, and percent of the
160 population as working age (18–64) as standalone metrics of socioeconomic vulnerability common
161 to many indices that have physical meaning. Note that while increasing index values signify
162 increasing vulnerability, increases in these standalone metrics signify decreasing vulnerability.

163 Infrastructural vulnerability, similar to exposure, includes metrics relevant to housing, such as
164 building age and construction type. Lastly, we include two metrics of flood experience, number of
165 federally declared disasters in the past three years and county-level participation in the Community
166 Rating System (CRS) program. The CRS program is a mechanism to incentivize insurance uptake
167 and increase community-level flood resilience through community-wide policy discounts offered
168 in exchange for flood risk reduction actions. 70% of NFIP policyholders nationwide live in
169 participating communities, and 26 out of California’s 58 counties have participated in the CRS at
170 some point since its inception, with 24 of those counties still in the program today (FEMA 2023a).

171 **Spatiotemporal Resolution**

172 To analyze the effects of spatial and temporal resolution on predictive accuracy and model
173 interpretability results, we fit RF models at two spatial and two temporal scales, for an overall
174 total of four models. The two spatial scales are at the county-level across all of California and at
175 the census tract-level across Sacramento County. Sacramento County was chosen because of its
176 history of flood events (James and Singer 2008) and its significant investment in flood mitigation,
177 particularly its commitment to the CRS program. The two temporal scales are 1981–2021 and
178 2009–2021. While most hazard variables are available starting in 1981, many exposure and
179 vulnerability variables are not available until later. For the models starting in 2009, the variables
180 footnoted with (b) or (c) in Table 1 no longer require extrapolation beyond the range of their record
181 and the variables footnoted with (d) are allowed to vary annually rather than remain constant.

182 Nearly all of the hazard variables have higher temporal resolutions than the exposure and
183 vulnerability variables, even in the 2009–2021 timeframe. Most hazard variables are recorded at
184 the event scale, meaning that each record in the dataset will have different values based on the

185 characteristics (maximum IVT, 3-day antecedent soil moisture, etc.) of that particular AR event.
186 The exposure and vulnerability variables, though, are recorded at the annual scale. This is logical
187 in a physical sense; for example, the total number of housing units does not change day-to-day like
188 the weather does. From a data perspective, though, this means that there is no intra-annual variation
189 to exploit for the RF model, and RF models are known to preferentially split on features with higher
190 variance (Strobl et al. 2007). The significantly higher resolution of the hazard data is analogous to
191 the disproportionate focus on hazard in the physical model space (Merz et al. 2010). Nevertheless,
192 comparison of the 1981–2021 and 2009–2021 models provides some insight into the usefulness of
193 collecting additional exposure and vulnerability data for the earlier years in the historical record.

194 **METHODOLOGY**

195 **Dataset Preparation**

196 We generated four datasets of AR events spanning the forty-year time period from 1981–2021.
197 Each record included the AR characteristics of the event plus the additional hazard, exposure, and
198 vulnerability variables documented for that specific time and place. If an AR passed over multiple
199 geographic units, it was tabulated as multiple records (one per county/tract) in order to capture the
200 effect of differences in exposure and vulnerability between different locations. Table 2 shows the
201 total number of records in each of the four datasets.

202 We sampled 80% of the data using stratified random sampling by county/tract for training and
203 validation, reserving the remaining 20% to test the performance of the final model. Stratified random
204 sampling ensures that there are counties (in the statewide model) or tracts (in the Sacramento model)
205 the model has never seen before, which in turn ensures that the goodness-of-fit metrics calculated
206 on the test set are more faithful representations of the model’s true performance. We then
207 implemented the Synthetic Minority Oversampling TEchnique (SMOTE) (Chawla et al. 2002) on
208 the training data to fix the class imbalance in Table 2. With highly imbalanced classes, it is difficult
209 to train an ML model that accurately captures damaging events; simply put, a naive model that
210 predicts no damage every time would be correct approximately 95% (statewide models) or 99.5%
211 (Sacramento models) of the time. SMOTE increases (oversamples) the number of records in the

212 minority class, in our case the damaging AR events, by creating synthetic records based on the
213 distribution of historical events, and decreases (undersamples) the number of non-damaging AR
214 events by a corresponding amount to achieve an even class balance. Combined over/undersampling
215 on imbalanced data improves the predictive accuracy of ML models across a range of contexts (e.g.,
216 [Estabrooks et al. \(2004\)](#)).

217 After implementing SMOTE on the training set for each model, we performed feature selection
218 to remove highly collinear variables. [While feature collinearity does not affect the accuracy of RF](#)
219 [models, it does adversely impact model interpretability, which is the main goal of this paper. We](#)
220 [therefore first identified clusters of correlated variables using principal component analysis \(PCA\),](#)
221 [then calculated the Akaike information criterion \(AIC\) of each variable in the cluster and kept only](#)
222 [the ones with the highest predictive power.](#) We repeated this [two-step](#) process of clustering and
223 consolidating until the maximum variable inflation factor (VIF) fell below 10 and the maximum
224 Pearson correlation coefficient fell below 0.8 ([James et al. 2013](#)). Despite the stochastic nature of
225 the SMOTE algorithm, the clusters were very stable, and our process removed a consistent subset
226 of the variables every time. Finally, we added one additional feature with uniform random noise,
227 which serves as a check for our feature importance and impact analyses; if a given feature is less
228 important than random noise, it is discarded.

229 **Model Training**

230 An RF model has three hyperparameters determined by the user: the number of trees in the
231 forest, the depth of each tree, and the number of predictors selected to fit each tree. The number of
232 trees in the forest was held constant at $n = 1,000$, [consistent with other RF applications in similar](#)
233 [contexts \(e.g., \[Alipour et al. \\(2020\\)\]\(#\)\)](#), and each tree was allowed to reach its maximum possible
234 depth (1 data point at each leaf). The number of predictors selected to fit each tree was tuned
235 between 1 and 10. We fit all models using 10-fold cross-validation, fitting the model on 90% of the
236 training data and calculating accuracy metrics on the remaining 10%, then repeating that process
237 across all the folds of the data. The best-fit model was chosen based on accuracy, which is the
238 number of correct predictions divided by the total number of predictions.

Performance Evaluation

We compared each of the fitted RF models against their respective withheld test sets. The test data had the same class imbalance as the original data, so we utilized three performance evaluation metrics appropriate for imbalanced data: area under the Receiver Operating Characteristic curve (ROC-AUC), area under the precision-recall curve (PR-AUC), and balanced accuracy. All three are derived from the confusion matrix, which summarizes the four potential outcomes for each prediction: true positive (TP), true negative (TN), false positive (FP), and false negative (FN). True positives and true negatives occur when the model correctly predicts a damaging or non-damaging AR event, respectively. False positives occur when the model incorrectly labels an AR as a damaging event, and false negatives occur when the model incorrectly labels an AR as a non-damaging event. From the confusion matrix, we derive secondary performance metrics, as shown in Equation 1. Precision measures correctly predicted positives out of all predicted positives. Recall, or sensitivity, measures correctly predicted positives out of all observed positives; the two names come from different disciplinary conventions, so we use both here in their respective contexts. Specificity measures correctly predicted negatives out of all observed negatives.

$$Precision = \frac{TP}{TP + FP}, \quad Recall \text{ (Sensitivity)} = \frac{TP}{TP + FN}, \quad Specificity = \frac{TN}{TN + FP}. \quad (1)$$

Sensitivity and specificity together are used to plot the ROC curve (Fig. 1a), a scale-invariant model diagnostic frequently used for imbalanced data (Kotsiantis et al. 2006). Precision and recall together are used to plot the PR curve (Fig. 1b). PR curves have been put forward as an even more informative tool for imbalanced datasets (Saito and Rehmsmeier 2015). The plots in Figure 1 measure the respective values of precision, recall (sensitivity), and specificity at varying detection thresholds for each of the four RF models. Model-optimized detection thresholds are marked with a dot. Finally, the balanced accuracy is defined as the average of sensitivity and specificity and calculated based on the model-optimized detection threshold. The values of the three metrics for each model are reported in Table 3. The fitted RFs at all spatial and temporal scales clearly

264 outperform the null models, signalling that they are far more informative about damaging events.

265 **MODEL INTERPRETATION**

266 **SHapley Additive exPlanations (SHAP)**

267 We use SHapley Additive exPlanations (SHAP; [Lundberg and Lee \(2017\)](#)), which utilize
268 Shapley values to assess local and global feature importance, for interpretation of our RF models.
269 Shapley values are a game-theory approach for fairly apportioning a “prize” between multiple
270 “players.” In the context of ML, the “players” are predictor variables and the “prize” is the difference
271 between the overall mean prediction (expected value) and the prediction for a specific observation
272 (observed value). SHAP conceptualizes the calculation of Shapley values as an additive feature
273 attribution model and estimates the feature’s contribution to the difference between the expected
274 and observed value for every record in the dataset. We prefer SHAP over other interpretable ML
275 techniques for several reasons. First, it is the only method to satisfy the three statistical properties
276 (local accuracy, missingness, and consistency) that are necessary and sufficient to ensure a fair
277 apportionment of the overall contribution among the various predictors ([Lundberg and Lee 2017](#)).
278 Second, it provides one coherent framework for examining both the magnitude and direction of
279 a feature’s effect and for building from local to global importance ([Molnar 2023](#)). Many other
280 interpretable ML methods only apply to one of these use cases, leading to an analysis that relies on
281 unrelated tools with different baseline assumptions. Third, comparisons between SHAP and other
282 interpretable ML techniques showed agreement at every step of the analysis.

283 **Feature Importance**

284 We first focus on global feature importance, which is calculated as the mean of the absolute value
285 of SHAP values across all observations. Figure 2 shows the overall importance of features in each
286 model, grouped by risk dimension and concept and normalized to a total of 100%. Across the four
287 models, hazard features account for approximately 70% of the model’s predictive power, exposure
288 features account for approximately 10%, and vulnerability features comprise the remaining 20%.
289 The statewide 1981–2021 model is an outlier with higher-than-average exposure and vulnerability

290 contributions. Figure 2 also separates the risk dimensions by the concepts defined in Table 1. We
291 notice some patterns; for example, social vulnerability is the most important vulnerability concept
292 in the statewide models, and flood experience is more important in the 1981 models than the 2009
293 models. Population exposure is the most important exposure concept in the Sacramento models,
294 while housing exposure is more important in the statewide models. Factors affecting insurance
295 take-up only appear as important in the 1981 models. The relative contributions of the hazard
296 concepts are more stable across the spatial and temporal scales, but climate modes and land surface
297 variables are slightly more important in the 2009 models.

298 We move to comparing the rankings of individual features, noting common trends across all
299 four models and exploring differences. In all cases, the top five predictors are related to the hazard
300 dimension, mostly from either the AR characteristics concept or the antecedent conditions concept.
301 Total precipitation and maximum IVT always occupy the top two positions, and total precipitation
302 is the most important predictor in three out of four models. Lagged cumulative precipitation and
303 lagged average soil moisture appear frequently, which highlights the important contribution of
304 antecedent conditions to AR-driven flood risk in California. Two features related to exposure are
305 noteworthy: percentage of the population living in the FEMA 100-year floodplain appears in the
306 top ten in all four models, and total number of housing units appears in the top ten in both statewide
307 models. For vulnerability, CRS score occupies the tenth position in the 1981 statewide model and
308 median housing unit age occupies the tenth position in the 1981 Sacramento model. There are no
309 features related to the vulnerability dimension in the top ten of either of the 2009 models.

310 There are more noticeable shifts in feature rank across space than across time. While it is
311 important in all cases, total precipitation has a larger SHAP feature contribution in the statewide
312 models. The contributions are more evenly spread among the top five to ten features in the
313 Sacramento models. AR category is more important in the 1981 models and ENSO climate index
314 is more important in the 2009 models, but otherwise there are no clear temporal patterns. While
315 our methodology does not necessarily allow for a direct comparison between the two, this may
316 suggest that spatial scale has a larger effect on model results than temporal scale.

317 **Feature Impact**

318 Feature impact plots visualize the direction and magnitude of a particular feature's influence
319 on the model's prediction. We use accumulated local effects (ALE) curves (Apley and Zhu 2020)
320 paired with SHAP values to analyze feature impact. The atomic unit of an ALE curve is the local
321 effect, or the difference between the prediction at x and the prediction at some perturbed value of x
322 within a small interval δ . ALEs are the sum of the local effects for all observations falling within
323 $x \pm \delta$, calculated for each x in the feature domain. ALE curves show marginal contributions, not
324 conditional contributions, so the effects of correlated features are not separated; however, they are
325 more robust to collinearity than the more commonly used partial dependence plots (Stein et al.
326 2021), so they are well suited for our analysis. We pair the ALE curves, which are average metrics,
327 with random samples of 1,000 SHAP values from individual records. The SHAP values provide
328 an understanding of the variance and show where the ALE curve is based on more or less data.

329 Figures 4a–c plot the SHAP values and ALE curves for the top three hazard features in the
330 1981 statewide model: total precipitation (4a), AR maximum IVT (4b), and AR duration (4c). The
331 behavior of all of the hazard features follows a similar pattern, where probability of damage increases
332 with increasing feature values until some threshold point. At a certain point, the magnitude of the
333 hazard becomes so large that damage becomes the probable outcome. For total precipitation (Fig.
334 4a), the threshold point is roughly 75mm, or approximately 15% of California's mean annual total
335 precipitation. The threshold point for maximum IVT (Fig. 4b) is about 750mm, which would be a
336 Category 2–3 AR event, and the threshold point for AR duration (Fig. 4c) is about 30 hours.

337 Figures 4d–f plot the SHAP values and ALE curves for the top three exposure features in
338 the 1981 statewide model: total number of housing units (4d), percent population living in the
339 FEMA 100-year floodplain (4e), and percent housing stock as single family homes (4f). The plot
340 of total housing units (Fig. 4d) largely shows increasing probability of damage with increased
341 housing stock. More people and more buildings at risk imply more chances for damage, so this
342 matches intuitive reasoning. Less intuitive is the influence of the percent population living in the
343 100-year floodplain (Fig. 4e). We would expect more people in the floodplain to increase the

344 likelihood of damage; instead, it appears to have a negative influence on damage probability. The
345 distributions of the SHAP values in these panels provide more information about the discrepancy.
346 Statewide, the median percent population in the floodplain is 8.1%, which means about half of all
347 counties fall on the portion of the ALE curve in Figure 4b that is increasing. There could also be a
348 confounding relationship with county-level flood resilience; the counties with the highest percentage
349 of population living in the floodplain, say 20% or more, may be more prepared for flooding and thus
350 less likely to sustain damage from an AR event. Another counterintuitive relationship is the negative
351 influence of the percent of single family homes (Fig. 4f). NFIP policyholders disproportionately
352 live in single family homes, so more of this housing type would mean more opportunity for claims.
353 But residents of single family homes also tend to be more likely to invest in individual-level flood
354 mitigation efforts, so the negative relationship might indicate that this particular feature is capturing
355 more vulnerability than exposure. Percent single family homes is also negatively correlated with
356 total population, so counties with higher percentages likely have fewer opportunities to submit
357 insurance claims that would be recorded as damage.

358 Figures 4g–i plot the SHAP values and ALE curves for the top three vulnerability features in
359 the 1981 statewide model: CRS score (4g), median household income (4h), and CalEnviroScreen
360 pollution burden score (4i). Lower numbers indicate more significant flood resilience investment.
361 1 is the best possible score score and 10 means that a county has not engaged with the CRS. The
362 perhaps-surprising trend of higher damage probabilities at CRS scores of 7–9 compared to those
363 at a CRS score of 10 may be because counties with a history of damaging flood events are more
364 likely to invest time and money into joining the CRS program. Multiple studies have found that
365 CRS-participating counties and communities see significant reductions in flood loss (Highfield and
366 Brody 2017; Gourevitch and Pinter 2023), and the SHAP values in Figure 4g suggest that achieving
367 a CRS score of 6 or better does pay off in terms of flood risk reduction; however, Sacramento
368 County is the only county in California that has achieved a rating of 4 or better, so scores beyond
369 this point are not necessarily representative of the entire state. The median household income (Fig.
370 4h) and the CalEnviroScreen pollution burden score (Fig. 4i), both measures of social vulnerability,

371 seem to indicate that the probability of damage decreases with increasing vulnerability. This may
372 be indicative of limitations in the link between NFIP claims and damage sustained by all members
373 of the population, suggesting that broad interpretations for flood risk should be made with caution.
374 Further considerations for researchers using NFIP data are included in the Discussion.

375 **Feature Interactions**

376 Figure 5 plots interactions between hazard, as measured by the [Ralph et al. \(2019\)](#) intensity
377 category, and exposure and vulnerability. [In most cases, exposure and vulnerability variables](#)
378 [become less useful predictors of damage likelihood as hazard intensity increases.](#) Category 5 ARs
379 [have historically almost always caused some amount of damage \(Corringham et al. 2019\),](#) so while
380 [exposure and vulnerability variables may still impact the severity of damage,](#) they no longer have
381 [any influence on the probability.](#) For example, in both the total number of housing units (Fig. 5a)
382 and the percentage of housing stock as single-family homes (Fig. 5c), the observed trend is strongest
383 for Category 1 events and weakest for Category 5 events, when probability of damage is elevated
384 and individual-level characteristics are more likely to be overwhelmed by the severity of the hazard.
385 The pattern flips, though, for CRS score (Fig. 5d) [and median household income \(Fig. 5e\).](#) The
386 increase in probability of damage moving from a CRS score of 10 to 9 all but disappears for the
387 largest ARs, and the benefit of improving a county's CRS score from a 5 to a 4 or better increases
388 with increasing intensity category. [For median household income, there is a slight reduction in](#)
389 [damage probability for the lowest incomes that only occurs at the highest hazard intensities.](#) These
390 figures do not provide a comprehensive list of the ways hazard, exposure, and vulnerability interact;
391 rather, they illustrate examples of the kind of practically relevant insights from our RF models that
392 can be used to help communities better understand their risk under different AR scenarios.

393 **DISCUSSION**

394 **Benefits of Data-Driven Approach**

395 Through a combination of global feature importance, feature impact plots, and feature inter-
396 actions, our RF models were able to identify new connections between the risk dimensions and

397 AR-driven flood damage in California. The global feature importance analysis showed that hazard
398 features, individually and collectively, had the biggest influence on flood damage. This is not
399 surprising; flooding and flood damage are far more likely to occur on rainy days than sunny days.
400 An example of decision-relevant information related to hazard from our models are the nonlinear
401 thresholds in the feature impact plots. If emergency managers are confident that flood damage is
402 likely above a certain precipitation threshold, it reduces the information burden required to make
403 a decision and allows for quicker mobilization of resources. Another example comes from the
404 interaction plots: changes in exposure and vulnerability only affected damage probability during
405 Category 5 events in one out of four features shown, so a well-rounded flood resilience strategy
406 would include elements that decrease risk across the spectrum of potential hazards.

407 We also showed that exposure and vulnerability explain up to a third of the model's predictions.
408 This finding has important implications from a management perspective because exposure and
409 vulnerability can be altered at the community level more easily than hazard. The ALE-SHAP
410 plots are therefore useful to understand the magnitude and direction of effects and to highlight
411 overlooked risk factors for NFIP policyholders. For example, for the percent of the population
412 living in the 100-year floodplain (Fig. 4e), probability of damage peaks at 5–10% before starting to
413 decrease. Values in this range might indicate potential risk hotspots, where hazard is high enough
414 to cause damage but not high enough that counties have significantly invested in resilience efforts,
415 and identify a subset of counties that are worthy of more in-depth local analysis.

416 **Geographic Representation**

417 While random forests and interpretable ML methods are powerful tools for combining disparate
418 data sources and extracting insights, the results of the models are only as good as the data used
419 to train, fit, and validate them. We discuss two key limitations of this work, which stem from
420 the assumption that NFIP claims accurately represent total losses from floods and relate to (a)
421 geographic and (b) socioeconomic differences in insurance takeup rates.

422 First, there are significant inter-county differences in NFIP takeup rates, ranging from almost
423 0% in Mariposa County to over 18% in Sutter County. We compare observed versus expected

424 takeup rate (see Appendix I) to identify counties with anomalously high or low numbers of NFIP
425 policies-in-force. Sutter, Yuba, and Sacramento Counties all have more policies than expected.
426 All three border the Sacramento River, which has an extensive history of severe floods, and all
427 three have invested heavily in the CRS program: Yuba County has a score of 6 (top 30% of all
428 participating communities), Sutter County has a score of 5 (top 13%), and Sacramento County
429 has a score of 2 (top 0.5%) (FEMA 2023a). One of the activities that garners CRS points is
430 public advertisement of the NFIP, and better community ratings lead to higher policy discounts,
431 so CRS participation has likely increased the number of policies in these areas. Joining the CRS
432 also requires significant upfront investment, so it is more widely adopted in the communities and
433 counties that can afford to participate (Sadiq et al. 2020). On the other hand, Alpine, Mariposa,
434 Tuolumne, and Imperial Counties are rural counties with relatively low populations, and all have
435 60 or fewer policies-in-force in 2021. The low number of policies mean that even if a flood event
436 does cause damage or loss, it is less likely to lead to a NFIP insurance claim and therefore less
437 likely to be labeled as damaging event in our dataset. Future work could use the expected takeup
438 rates calculated here and determine the appropriate county-level correction factors to account for
439 differences in takeup rates that cannot be attributed to real differences in flood hazard.

440 **Socioeconomic Representation**

441 In addition to the unequal representation of different counties, the demographic and socioe-
442 conomic characteristics of NFIP policyholders contribute to representativeness issues. The most
443 socioeconomically vulnerable populations are typically most affected by floods and other disasters
444 (e.g., Debbage (2019)), but there are several intersectional factors affecting vulnerable populations
445 that simultaneously reduce the likelihood of NFIP participation and exacerbate flood risk. As
446 one example, renters are particularly vulnerable to negative consequences from flooding (Heiman
447 2022). People of color are more likely to be renters, and renters tend to have lower incomes than
448 homeowners (ACS 2023b). However, 80% of NFIP policyholders are single-family homeowners.
449 While the NFIP does offer policies for renters, renters are often unaware that standard renter's
450 insurance does not cover flood damage, and in many places landlords are not required to disclose

451 an apartment's history of flooding (Heiman 2022). Therefore renters are likely underrepresented
452 in NFIP claims data, meaning that if there are factors specific to renters that affect flood risk, they
453 will not be identified as important by our model. Using a different dataset than the NFIP as the
454 response variable, such as remote sensing imagery (Szczyrba et al. 2021) or post-event surveys
455 (Merz et al. 2013), could improve the representation of renters and other vulnerable populations.
456 The predictor variables used in this study could also be enhanced through interviews or community-
457 specific knowledge to better capture unaccounted-for resilience characteristics (Ismail-Zadeh et al.
458 2017). Lastly, using ML techniques to move from county- or state-level summary statistics to maps
459 of spatially varying hazard, exposure, and vulnerability would be of great benefit for future flood
460 mitigation investment decisions.

461 CONCLUSION

462 In this paper, we used interpretable machine learning (ML) tools to understand how the three
463 dimensions of risk, hazard, exposure, and vulnerability, relate to AR-induced flood damage in
464 California. We collected a large dataset of over forty predictor variables to quantify the contributions
465 of each the three dimensions to the probability of flood damage, as measured using flood insurance
466 claims from the National Flood Insurance Program (NFIP). We considered two spatial resolutions
467 for the data: the county scale, modeled for all of California, and the census tract scale, modeled for
468 Sacramento County. We also considered timescales of 1981–2021, using exposure and vulnerability
469 data with limited temporal variation, and 2009–2021, using exposure and vulnerability data at an
470 annual resolution. This produced a total of four random forest classification models, each of
471 which detected true positives (AR events with NFIP claims) with a high level of accuracy in very
472 imbalanced datasets.

473 We showed the power of interpretable ML to identify and investigate drivers of AR-driven
474 flood risk given publicly available hazard, exposure, and vulnerability data. We gained insight
475 into damage drivers by examining feature importance (how much does a given feature influence
476 the model's predictions?) and feature impact (how does increasing or decreasing the value of the
477 feature affect the response?) using SHapley Additive exPlanations (SHAP) as a unifying framework.

478 While hazard intensity features were the most important predictors of whether an AR would cause
479 damage, exposure and vulnerability contributed up to a third of the model's explanatory power,
480 and the overall relative contributions by risk dimension and risk concept broadly agreed across the
481 spatial and temporal scales considered. Total precipitation was the most important predictor in
482 three out of four models, and features related to the intensity of the hazard consistently represented
483 the majority of the top ten. An analysis of feature impact for the top three hazard features in
484 the county-level (statewide) model fit on data from 1981 onward revealed that increasing hazard
485 severity increased the probability of flood damage, up to some threshold point. Above that threshold,
486 probability of damage reached a saturation point where it was no longer sensitive to changes in
487 precipitation; at a certain point, it became not a question of if damage will occur, but how much.
488 In most cases, increasing exposure and vulnerability also increased the probability of damage,
489 although the interpretation differed slightly depending on the specific feature under consideration.
490 The physically plausible explanations of the data-driven outputs from SHAP and other interpretable
491 ML tools support our confidence that the model is characterizing real drivers of flood risk.

492 We also illustrated the ramifications of the assumptions made to fit our RF models utilizing
493 available data. Changes in the spatial and temporal resolution of the input data altered the ranking
494 of which features were deemed significant in the analysis of global feature importance. The higher
495 temporal scale of hazard data relative to the other risk dimensions and the differences in NFIP
496 representativeness across geographic and socioeconomic boundaries were noted as limitations. We
497 proposed avenues for future work that would mitigate these limitations and potentially uncover new
498 pathways to increased resilience. Overall, our work highlights both the possibilities and pitfalls of
499 using interpretable ML for flood risk assessment. It enhances our understanding of the relation-
500 ship between individual AR events and their negative effects and broadens the discussion around
501 AR-driven flood damage to include more explicit characterization of exposure and vulnerability.
502 Understanding the drivers of damage improves our ability to predict and prepare for the impacts of
503 ARs, today and in the future.

504 **Data Availability Statement**

505 All data used in this study is publicly available, and all code created to generate results is
506 available in a Github repository (Bowers 2023). In particular, the datasets described in Table 1
507 are available as downloadable CSV files, the *reproduce_figures.html* markdown file recreates all
508 figures and numerical results from this paper, and the *figure4.html* markdown file recreates Figure
509 4 for models at all spatial and temporal resolutions.

510 **Acknowledgments**

511 This material is based upon work supported by both the Stanford Graduate Fellowship and the
512 National Science Foundation (NSF) Graduate Research Fellowship under Grant No. 1000265549.
513 Any opinions, findings, and conclusions expressed in this material are those of the authors and
514 do not necessarily reflect the views of the NSF. CB contributed conceptualization, data curation,
515 methodology, formal analysis, validation, visualization, writing - original draft, and writing - review
516 and editing. KAS and JWB contributed supervision and writing - review and editing, and JWB
517 provided resources. We additionally thank Jenny Suckale and two anonymous reviewers for their
518 helpful feedback that improved the quality of this work.

APPENDIX I. CALCULATION OF NFIP TAKEUP RATE BY COUNTY

We calculated the observed and expected number of NFIP policies and insurance takeup rates for each county in California to determine which counties had more or fewer policies than predicted. Observed numbers of policies were calculated based on 2021 policies-in-force (FEMA 2023b), and observed takeup rates were calculated as the number of policies divided by the number of 2021 housing units in each county.

The expected numbers of policies and takeup rates were calculated as follows. We categorized both NFIP policies and housing units as either within-floodplain or out-of-floodplain based on the FEMA National Flood Hazard Layer (NFHL). NFIP policies were determined to be in or out of the floodplain by the NFHL flood zone code included in the policy information. Housing units were determined to be in or out of the floodplain by finding the percentage of each census block group that overlapped with a NFHL spatial polygon, then dividing the housing units in that block group assuming an even distribution in space. For example, if 40% of a block group was covered by the NFHL, then 40% of the housing units were labeled as in-floodplain (HU_{in}) and 60% were labeled as out-of-floodplain (HU_{out}). We summed all policies and housing units to estimate statewide within-floodplain and out-of-floodplain insurance takeup rates. The 2021 statewide within-floodplain takeup rate was found to be 12.6% and the 2021 statewide out-of-floodplain takeup rate was found to be 0.69%. The expected number of policies by county was then calculated by aggregating over all block groups in that county, as illustrated in Equation 2. Lastly, county-level expected takeup rates were found by dividing the expected number of policies by the total number of housing units in each county.

$$\text{Expected Policies} = \sum_{bg \in \left\{ \begin{array}{l} \text{all block} \\ \text{groups} \end{array} \right\}} 0.126 * (HU_{in})_{bg} + 0.0069 * (HU_{out})_{bg} \quad (2)$$

REFERENCES

- ACS (2023a). “DP04: Selected Housing Characteristics, 2009-2021, <<https://data.census.gov/>>.
- ACS (2023b). “DP05: ACS Demographic and Housing Estimates, 2009-2021, <<https://data.census.gov/>>.
- Alipour, A., Ahmadalipour, A., Abbaszadeh, P., and Moradkhani, H. (2020). “Leveraging machine learning for predicting flash flood damage in the Southeast US.” *Environmental Research Letters*, 15(2), 024011.
- Apley, D. W. and Zhu, J. (2020). “Visualizing the Effects of Predictor Variables in Black Box Supervised Learning Models.” *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 82(4), 1059–1086.
- Atreya, A., Ferreira, S., and Michel-Kerjan, E. (2015). “What drives households to buy flood insurance? New evidence from Georgia.” *Ecological Economics*, 117, 153–161.
- August, L., Bangia, K., Plummer, L., Prasad, S., Ranjbar, K., Slocombe, A., and Wieland, W. (2021). “CalEnviroScreen 4.0.” *Report no.*, California Office of Environmental Health Hazard Assessment, Sacramento, CA, <<https://oehha.ca.gov/media/downloads/calenviroscreen/report/calenviroscreen40reportf2021.pdf>>.
- Bakkensen, L. A., Fox-Lent, C., Read, L. K., and Linkov, I. (2017). “Validating Resilience and Vulnerability Indices in the Context of Natural Disasters.” *Risk Analysis*, 37(5), 982–1004.
- Bergstrand, K., Mayer, B., Brumback, B., and Zhang, Y. (2015). “Assessing the Relationship Between Social Vulnerability and Community Resilience to Hazards.” *Social Indicators Research*, 122(2), 391–409.
- Blum, A. G., Ferraro, P. J., Archfield, S. A., and Ryberg, K. R. (2020). “Causal Effect of Impervious Cover on Annual Flood Magnitude for the United States.” *Geophysical Research Letters*, 47(5).
- Bowers, C. (2023). “Supplemental Code Release: Uncovering Effects of Exposure and Vulnerability on Atmospheric River Flood Damage using Interpretable Machine Learning, <<https://github.com/corinnebowers/damagedrivers>>.
- Bradt, J. T., Kousky, C., and Wing, O. E. (2021). “Voluntary purchases and adverse selection in

568 the market for flood insurance.” *Journal of Environmental Economics and Management*, 110,
569 102515.

570 Brody, S. D., Kim, H., and Gunn, J. (2013). “Examining the Impacts of Development Patterns on
571 Flooding on the Gulf of Mexico Coast.” *Urban Studies*, 50(4), 789–806.

572 Brunner, M. I., Sikorska, A. E., and Seibert, J. (2018). “Bivariate analysis of floods in climate
573 impact assessments.” *Science of The Total Environment*, 616-617, 1392–1403.

574 Cao, Q., Gershunov, A., Shulgina, T., Ralph, F. M., Sun, N., and Lettenmaier, D. P. (2020). “Floods
575 due to Atmospheric Rivers along the U.S. West Coast: The Role of Antecedent Soil Moisture in
576 a Warming Climate.” *Journal of Hydrometeorology*, 21(8), 1827–1845.

577 CDC (2022). “CDC/ATSDR Social Vulnerability Index, 2000-2020,
578 <<https://www.atsdr.cdc.gov/placeandhealth/svi/data>

579 Chawla, N. V., Bowyer, K. W., Hall, L. O., and Kegelmeyer, W. P. (2002). “SMOTE: Synthetic
580 Minority Over-sampling Technique.” *Journal of Artificial Intelligence Research*, 16, 321–357.

581 Corringham, T. W. and Cayan, D. R. (2019). “The Effect of El Niño on Flood Damages in the
582 Western United States.” *Weather, Climate, and Society*, 11(3), 489–504.

583 Corringham, T. W., Ralph, F. M., Gershunov, A., Cayan, D. R., and Talbot, C. A. (2019). “Atmo-
584 spheric rivers drive flood damages in the western United States.” *Science Advances*, 5(12).

585 Cutter, S. L. (2016). “The landscape of disaster resilience indicators in the USA.” *Natural Hazards*,
586 80(2), 741–758.

587 Czajkowski, J., Villarini, G., Montgomery, M., Michel-Kerjan, E., and Goska, R. (2017). “As-
588 ssuming Current and Future Freshwater Flood Risk from North Atlantic Tropical Cyclones via
589 Insurance Claims.” *Scientific Reports*, 7(1), 41609.

590 Darlington, J. C. and Yiannakoulis, N. (2022). “Experimental Evidence for Coverage Preferences
591 in Flood Insurance.” *International Journal of Disaster Risk Science*, 13(2), 178–189.

592 Debbage, N. (2019). “Multiscalar spatial analysis of urban flood risk and environmental justice in
593 the Charlanta megaregion, USA.” *Anthropocene*, 28, 100226.

594 DeFlorio, M. J., Pierce, D. W., Cayan, D. R., and Miller, A. J. (2013). “Western U.S. extreme

595 precipitation events and their relation to ENSO and PDO in CCSM4.” *Journal of Climate*,
596 26(12), 4231–4243.

597 Dewitz, J. and USGS (2021). “National Land Cover Database (NLCD) 2019 Products (ver. 2.0,
598 June 2021).

599 Erlingis, J. M., Rodell, M., Peters-Lidard, C. D., Li, B., Kumar, S. V., Famiglietti, J. S., Granger,
600 S. L., Hurley, J. V., Liu, P., and Mocko, D. M. (2021). “A High-Resolution Land Data Assimilation
601 System Optimized for the Western United States [Dataset].” *Journal of the American Water
602 Resources Association*, 57(5), 692–710.

603 Estabrooks, A., Jo, T., and Japkowicz, N. (2004). “A multiple resampling method for learning from
604 imbalanced data sets.” *Computational Intelligence*, 20(1), 18–36.

605 FEMA (2006). “Hazus Flood Model Technical Manual.” *Report no.*, Department of Home-
606 land Security, Washington, DC, <[https://www.fema.gov/flood-maps/tools-resources/flood-map-
607 products/hazus/user-technical-manuals](https://www.fema.gov/flood-maps/tools-resources/flood-map-products/hazus/user-technical-manuals)>.

608 FEMA (2020). “National Flood Hazard Layer (NFHL), <[https://www.fema.gov/flood-maps/tools-
609 resources/flood-map-products/national-flood-hazard-layer](https://www.fema.gov/flood-maps/tools-resources/flood-map-products/national-flood-hazard-layer)>.

610 FEMA (2023a). “CRS Participating Communities, <[https://www.fema.gov/floodplain-
611 management/community-rating-system](https://www.fema.gov/floodplain-management/community-rating-system)>.

612 FEMA (2023b). “OpenFEMA Data Sets, <<https://www.fema.gov/about/openfema/data-sets>>.

613 Flanagan, B. E., Gregory, E. W., Hallisey, E. J., Heitgerd, J. L., and Lewis, B. (2011). “A Social
614 Vulnerability Index for Disaster Management.” *Journal of Homeland Security and Emergency
615 Management*, 8(1).

616 Gelaro, R., McCarty, W., Suárez, M. J., Todling, R., Molod, A., Takacs, L., Randles, C. A.,
617 Darmenov, A., Bosilovich, M. G., Reichle, R., Wargan, K., Coy, L., Cullather, R., Draper, C.,
618 Akella, S., Buchard, V., Conaty, A., da Silva, A. M., Gu, W., Kim, G.-K., Koster, R., Lucchesi, R.,
619 Merkova, D., Nielsen, J. E., Partyka, G., Pawson, S., Putman, W., Rienecker, M., Schubert, S. D.,
620 Sienkiewicz, M., and Zhao, B. (2017). “The Modern-Era Retrospective Analysis for Research
621 and Applications, Version 2 (MERRA-2) [Dataset].” *Journal of Climate*, 30(14), 5419–5454.

622 Gourevitch, J. D. and Pinter, N. (2023). “Federal incentives for community-level climate adaptation:
623 an evaluation of FEMA’s Community Rating System.” *Environmental Research Letters*, 18(3).

624 Heiman, E. R. (2022). “Protecting Renters from Flood Loss.” *University of Pennsylvania Law*
625 *Review*, 170(3), 783–809.

626 Highfield, W. E. and Brody, S. D. (2017). “Determining the effects of the FEMA Community
627 Rating System program on flood losses in the United States.” *International Journal of Disaster*
628 *Risk Reduction*, 21(November 2016), 396–404.

629 Highfield, W. E., Brody, S. D., and Shepard, C. (2018). “The effects of estuarine wetlands on flood
630 losses associated with storm surge.” *Ocean & Coastal Management*, 157, 50–55.

631 Highfield, W. E., Peacock, W. G., and Van Zandt, S. (2014). “Mitigation Planning: Why Haz-
632 ard Exposure, Structural Vulnerability, and Social Vulnerability Matter.” *Journal of Planning*
633 *Education and Research*, 34(3), 287–300.

634 Ismail-Zadeh, A. T., Cutter, S. L., Takeuchi, K., and Paton, D. (2017). “Forging a paradigm shift in
635 disaster science.” *Natural Hazards*, 86, 969–988.

636 James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013). *An Introduction to Statistical Learning*.
637 Springer Texts in Statistics. Springer, New York, NY.

638 James, L. A. and Singer, M. B. (2008). “Development of the Lower Sacramento Valley Flood-
639 Control System: Historical Perspective.” *Natural Hazards Review*, 9(3), 125–135.

640 Jonkman, S. N. (2005). “Global perspectives on loss of human life caused by floods.” *Natural*
641 *Hazards*, 34(2), 151–175.

642 Knighton, J., Buchanan, B., Guzman, C., Elliott, R., White, E., and Rahm, B. (2020). “Predicting
643 flood insurance claims with hydrologic and socioeconomic demographics via machine learning:
644 Exploring the roles of topography, minority populations, and political dissimilarity.” *Journal of*
645 *Environmental Management*, 272, 111051.

646 Komolafe, A., Herath, S., and Avtar, R. (2018). “Sensitivity of flood damage estimation to spatial
647 resolution.” *Journal of Flood Risk Management*, 11, S370–S381.

648 Konrad, C. P. and Dettinger, M. D. (2017). “Flood Runoff in Relation to Water Vapor Transport by

649 Atmospheric Rivers Over the Western United States, 1949–2015.” *Geophysical Research Letters*,
650 44(22), 456–11.

651 Kotsiantis, S., Kanellopoulos, D., and Pintelas, P. (2006). “Handling imbalanced datasets: A
652 review.” *GESTS International Transactions on Computer Science and Engineering*, 30.

653 Kousky, C. (2011). “Understanding the Demand for Flood Insurance.” *Natural Hazards Review*,
654 12(2), 96–110.

655 Lamjiri, M. A., Dettinger, M. D., Ralph, F. M., and Guan, B. (2017). “Hourly storm characteristics
656 along the U.S. West Coast: Role of atmospheric rivers in extreme precipitation.” *Geophysical
657 Research Letters*, 44(13), 7020–7028.

658 Lundberg, S. M. and Lee, S.-i. (2017). “A Unified Approach to Interpreting Model Predictions.”
659 *31st Conference on Neural Information Processing Systems (NIPS 2017)*, Long Beach, CA.

660 Merz, B., Kreibich, H., and Lall, U. (2013). “Multi-variate flood damage assessment: a tree-based
661 data-mining approach.” *Natural Hazards and Earth System Sciences*, 13, 53–64.

662 Merz, B., Kreibich, H., Schwarze, R., and Thielen, A. H. (2010). “Review article: “Assessment of
663 economic flood damage”.” *Natural Hazards and Earth System Sciences*, 10(8), 1697–1724.

664 Mobley, W., Sebastian, A., Blessing, R., Highfield, W. E., Stearns, L., and Brody, S. D. (2021).
665 “Quantification of continuous flood hazard using random forest classification and flood insurance
666 claims at large spatial scales: a pilot study in southeast Texas.” *Natural Hazards and Earth System
667 Sciences*, 21(2), 807–822.

668 Molnar, C. (2023). “Interpretable Machine Learning: A Guide for Making Black Box Models
669 Explainable, <<https://christophm.github.io/interpretable-ml-book/>>.

670 NOAA. “Multivariate ENSO Index Version 2 (MEI.v2), <<https://psl.noaa.gov/enso/mei/>>.

671 NOAA. “Pacific Decadal Oscillation (PDO), <<https://www.ncei.noaa.gov/access/monitoring/pdo/>>.

672 NOAA NCEI (2023). “U.S. Billion-Dollar Weather and Climate Disasters,
673 <<https://www.ncei.noaa.gov/access/billions/>>.

674 PEP (2023). “Population and Housing Unit Estimates Tables, <[https://www.census.gov/programs-
675 surveys/popest/data/tables.html](https://www.census.gov/programs-surveys/popest/data/tables.html)>.

676 Pollack, A. B., Sue Wing, I., and Nolte, C. (2022). “Aggregation bias and its drivers in large-scale
677 flood loss estimation: A Massachusetts case study.” *Journal of Flood Risk Management*, 15(4),
678 1–16.

679 Ralph, F. M., Rutz, J. J., Cordeira, J. M., Dettinger, M. D., Anderson, M., Reynolds, D., Schick,
680 L. J., and Smallcomb, C. (2019). “A scale to characterize the strength and impacts of atmospheric
681 rivers.” *Bulletin of the American Meteorological Society*, 100(2), 269–289.

682 Rözer, V., Kreibich, H., Schröter, K., Müller, M., Sairam, N., Doss-Gollin, J., Lall, U., and Merz,
683 B. (2019). “Probabilistic Models Significantly Reduce Uncertainty in Hurricane Harvey Pluvial
684 Flood Loss Estimates.” *Earth’s Future*, 7(4), 384–394.

685 Rufat, S., Tate, E., Burton, C. G., and Maroof, A. S. (2015). “Social vulnerability to floods:
686 Review of case studies and implications for measurement.” *International Journal of Disaster
687 Risk Reduction*, 14, 470–486.

688 Rutz, J. J., Steenburgh, W. J., and Ralph, F. M. (2014). “Climatological characteristics of atmo-
689 spheric rivers and their inland penetration over the western united states.” *Monthly Weather
690 Review*, 142(2), 905–921.

691 Sadiq, A. A., Tyler, J., and Noonan, D. (2020). “Participation and non-participation in FEMA’s
692 Community Rating System (CRS) program: Insights from CRS coordinators and floodplain
693 managers.” *International Journal of Disaster Risk Reduction*, 48(September 2019).

694 Sadler, J., Goodall, J., Morsy, M., and Spencer, K. (2018). “Modeling urban coastal flood severity
695 from crowd-sourced flood reports using Poisson regression and Random Forest.” *Journal of
696 Hydrology*, 559, 43–55.

697 Saito, T. and Rehmsmeier, M. (2015). “The precision-recall plot is more informative than the ROC
698 plot when evaluating binary classifiers on imbalanced datasets.” *PLoS ONE*, 10(3), 1–21.

699 Sanders, B. F., Schubert, J. E., Kahl, D. T., Mach, K. J., Brady, D., AghaKouchak, A., Forman, F.,
700 Matthew, R. A., Ulibarri, N., and Davis, S. J. (2022). “Large and inequitable flood risks in Los
701 Angeles, California.” *Nature Sustainability*, 6(1), 47–57.

702 Smith, A. B. and Katz, R. W. (2013). “US billion-dollar weather and climate disasters: Data

703 sources, trends, accuracy and biases.” *Natural Hazards*, 67(2), 387–410.

704 Solomatine, D. P. and Ostfeld, A. (2008). “Data-driven modelling: Some past experiences and new
705 approaches.” *Journal of Hydroinformatics*, 10(1), 3–22.

706 Stein, L., Clark, M. P., Knoben, W. J., Pianosi, F., and Woods, R. A. (2021). “How Do Climate and
707 Catchment Attributes Influence Flood Generating Processes? A Large-Sample Study for 671
708 Catchments Across the Contiguous USA.” *Water Resources Research*, 57(4), 1–21.

709 Strobl, C., Boulesteix, A. L., Zeileis, A., and Hothorn, T. (2007). “Bias in random forest variable
710 importance measures: Illustrations, sources and a solution.” *BMC Bioinformatics*, 8.

711 Szczyrba, L., Zhang, Y., Pamukcu, D., Eroglu, D. I., and Weiss, R. (2021). “Quantifying the Role
712 of Vulnerability in Hurricane Damage via a Machine Learning Case Study.” *Natural Hazards
713 Review*, 22(3), 1–12.

714 Tate, E., Rahman, M. A., Emrich, C. T., and Sampson, C. C. (2021). “Flood exposure and social
715 vulnerability in the United States.” *Natural Hazards*, 106(1), 435–457.

716 Thomson, H., Zeff, H. B., Kleiman, R., Sebastian, A., and Characklis, G. W. (2023). “Systemic
717 Financial Risk Arising From Residential Flood Losses.” *Earth’s Future*, 11(4), 86.

718 USGS (2023). “National Hydrography Dataset, <[https://www.usgs.gov/national-
719 hydrography/access-national-hydrography-products](https://www.usgs.gov/national-hydrography/access-national-hydrography-products)>.

720 Wagenaar, D., de Jong, J., and Bouwer, L. M. (2017). “Multi-variable flood damage modelling with
721 limited data using supervised learning approaches.” *Natural Hazards and Earth System Sciences*,
722 17(9), 1683–1696.

723 Willis, H., Narayanan, A., Fischbach, J., Molina-Perez, E., Stelzner, C., Loa, K., and Kendrick, L.
724 (2016). “Current and Future Exposure of Infrastructure in the United States to Natural Hazards.”
725 *Report no.*, Santa Monica, CA.

726 Woldemeskel, F. and Sharma, A. (2016). “Should flood regimes change in a warming climate? The
727 role of antecedent moisture conditions.” *Geophysical Research Letters*, 43(14), 7556–7563.

728 Zahran, S., Weiler, S., Brody, S. D., Lindell, M. K., and Highfield, W. E. (2009). “Modeling
729 national flood insurance policy holding at the county scale in Florida, 1999–2005.” *Ecological*

Fig. 1. Model performance metrics for the four RF models. The spatial scale is represented by color, where Sacramento is blue and statewide is gold; the temporal scale is represented by shading, where darker indicates 1981–2021 and lighter indicates 2009–2021. **(a)** Receiver Operating Characteristic (ROC) curves. A perfect model would reach the top-left corner of the plot. The solid lines are the results for the respective models at different detection thresholds between 0 and 1, and the points on the curves indicate the sensitivity and specificity of the model-optimized detection threshold. The dashed black line represents the ROC of random guessing. **(b)** Precision-recall (PR) curves. A perfect model would reach the top-right corner of the plot. The solid lines are the results for the respective models at different detection thresholds between 0 and 1, and the points on the curves indicate the precision and recall of the model-optimized detection threshold. The dashed lines represent the precision of the null models that predict no damage for all records. The precision is constant and equal to the class imbalance ratio of the test data.

Fig. 2. SHAP relative global feature importance by risk dimension and concept. Relative global feature importance is shown for both spatial resolutions (statewide and Sacramento County) and both temporal resolutions (1981–2021 and 2009–2021). The color represents the risk dimension, and the bolded percentages indicate the overall contribution of that risk dimension to the overall model performance. The shading, labeled on the right-hand side of the plot, represents the risk concept as defined in Table 1. Feature-level SHAP importance estimates are normalized so that the total for each model sums to 100%.

Fig. 3. SHAP feature importance rank. Features are colored by their global SHAP feature contribution in each model, where the global SHAP feature contribution is calculated as the mean of all observation-level SHAP values. For example, a SHAP feature contribution of 5% means that the value of that feature increases or decreases the probability of damage by 5% on average. The top ten features with the largest SHAP feature contributions in each model are labeled with their rank. Features without shaded bars were either removed during feature selection or had smaller global SHAP values than the random noise variable. The H/E/V labels along the left side of the plot indicate whether each feature is related to the hazard, exposure, or vulnerability dimension.

Fig. 4. Top three hazard, exposure, and vulnerability features in the 1981 statewide model.

The top three most important hazard features are **(a)** total precipitation, **(b)** AR maximum IVT, and **(c)** AR duration. The top three most important exposure features are **(d)** total housing units, **(e)** percent of the population living in the 100-year floodplain, and **(f)** percent of the housing stock as single family homes. The top three most important vulnerability features are **(g)** CRS score, **(h)** median household income, and **(i)** CalEnviroScreen pollution burden score. The left and right Y-axes represent the change in damage probability relative to the average probability of damage. Black lines represent ALE curves that show the average trend between predictor and response; ALE curves are plotted over the middle 95% of the data to reduce the visual impact of outliers. Gray points represent SHAP values for 1,000 individual observations randomly sampled from the training set.

Fig. 5. Interactions between AR category and exposure and vulnerability features in the 1981 statewide model. Exposure features are (a) total housing units, (b) percent of the population living in the 100-year floodplain, and (c) percent of housing stock as single-family homes. Vulnerability features are (d) CalEnviroScreen Pollution Burden score, (e) median household income, and (f) CRS score. ARs of different categories are colored based on the legend at the bottom. The tick marks along the bottom of each panel indicate the distribution of values for that feature. ALE curves are plotted over 95% of the data to reduce the impact of outliers, and lines are plotted with a smoothing factor to reduce visual clutter.

TABLE 1. Predictor variables. This table includes variables that were retained through feature selection in at least one of the four models. Variables are grouped by risk dimension and variable concept. Spatial resolution is at the county scale (C), census tract scale (T), or constant (–). Temporal resolution is by event (E), monthly (M), annual (A), or constant (–). The Data Source column lists a citation where data for the variable can be retrieved, the Justification column lists a citation supporting the variable’s inclusion in the model.

Risk Dimension	Concept	Variable	Spatial	Temporal	Data Source	Justification
Hazard	AR characteristics	Maximum IVT (kg/m/s)	T ^a	E	Gelaro et al. (2017)	Corringham et al. (2019)
		Duration (hours)	T ^a	E	Rutz et al. (2014)	Corringham et al. (2019)
		AR category	T ^a	E	Ralph et al. (2019)	Corringham et al. (2019)
		Total precipitation (mm)	T ^a	E	Gelaro et al. (2017)	Brunner et al. (2018)
	An- tecedent conditions	3- & 14-day total precipitation prior to AR event (mm)	T ^a	E	Gelaro et al. (2017)	Woldemeskel and Sharma (2016)
		3-day mean soil moisture prior to AR event (mm/m)	T	E	Erlingis et al. (2021)	Cao et al. (2020)
	Climate modes	El Niño Southern Oscillation (ENSO)	–	M	NOAA (a)	Corringham and Cayan (2019)
		Pacific Decadal Oscillation (PDO)	–	M	NOAA (b)	DeFlorio et al. (2013)
	Land surface	Impervious land cover (%)	T	Y ^c	Dewitz and USGS (2021)	Blum et al. (2020)
		Developed land cover (%)	T	Y ^c	Dewitz and USGS (2021)	Brody et al. (2013)
Wetlands land cover (%)		T	Y ^c	Dewitz and USGS (2021)	Highfield et al. (2018)	
Exposure	Population	Population density per square mile	T ^b	Y	PEP (2023)	Jonkman (2005)
		Population within FEMA floodplain (%)	T	–	FEMA (2020)	Sanders et al. (2022)
	Housing	Housing units	T ^b	Y	PEP (2023)	Willis et al. (2016)
		Single-family homes (%)	T	– ^d	ACS (2023a)	Rufat et al. (2015)

Risk Dimension	Concept	Variable	Spatial	Temporal	Data Source	Justification
	Insurance take-up	Riverside census tract indicator	T	–	USGS (2023)	Kousky (2011)
		Coastal county indicator	C	–	USGS (2023)	Kousky (2011)
Vulnerability	Socioeconomic	CDC SVI components	T	Y ^c	CDC (2022)	Bakkensen et al. (2017)
		CalEnviroScreen metrics	T	–	August et al. (2021)	Bergstrand et al. (2015)
		Median household income (2022 dollars)	T ^b	Y	ACS (2023a)	Cutter (2016)
		Non-Hispanic white population (%)	T ^b	Y	ACS (2023b)	Cutter (2016)
		Working-age population (%)	T ^b	Y	ACS (2023b)	Cutter (2016)
	Infrastructural	Housing units over 40 years old (%)	T	– ^d	ACS (2023a)	Highfield et al. (2014)
		Median housing age (years)	T	– ^d	ACS (2023a)	Knighton et al. (2020)
		Owner-occupied housing units (%)	T	– ^d	ACS (2023a)	Rufat et al. (2015)
		Mobile homes (%)	T	– ^d	ACS (2023a)	Tate et al. (2021)
	Flood experience	3-year lagged total disaster declarations	C	Y	FEMA (2023b)	Bakkensen et al. (2017)
		Community Rating System (CRS) score	C	Y	FEMA (2023a)	Highfield and Brody (2017)

^a Values are aggregated to the tract/county level from MERRA-2 (~50×50 km grid cells) (Gelaro et al. 2017).

^b Pre-2000 tract-level estimates are based on a weighted distribution of county/statewide values, where weights are determined as a function of the distribution of tract values in 2000.

^c Data is only available starting from 2000 (for the CDC SVI) or 2001 (for the land surface variables), so for prior years the values are equal to those recorded in the earliest year of data.

^d Data is only available starting from 2009, so these are considered to be temporally constant in the 1981–2021 models and annually varying in the 2009–2021 models.

TABLE 2. Dataset statistics by spatial and temporal scale.

Spatial Scale	Temporal Scale	Number of Records	Damaging Records	Class Balance
Statewide	1981–2021	17,265	863	5.0%
Statewide	2009–2021	5,659	234	4.1%
Sacramento	1981–2021	145,977	874	0.60%
Sacramento	2009–2021	47,776	114	0.24%

TABLE 3. Random forest model performance metrics. Fitted model performance (*RF*) is compared against null model performance (*Null*) for each model and each metric under consideration.

Spatial Resolution	Temporal Resolution	ROC-AUC		PR-AUC		Balanced Accuracy	
		<i>RF</i>	<i>Null</i>	<i>RF</i>	<i>Null</i>	<i>RF</i>	<i>Null</i>
Statewide	1981–2021	0.914	0.500	0.435	0.051	0.707	0.500
Statewide	2009–2021	0.914	0.500	0.353	0.041	0.658	0.500
Sacramento	1981–2021	0.959	0.500	0.311	0.007	0.761	0.500
Sacramento	2009–2021	0.898	0.500	0.046	0.002	0.702	0.500