



## Abstract

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

1: Determining the harvest location of timber is crucial to enforcing international regulations designed to tackle illegal logging and associated trade in forest products. However, complex supply chains obscure harvest sources, which often leaves paper-based traceability systems as the sole tool for demonstrating provenance, despite its vulnerability to fraud. Stable Isotope Ratio Analysis (SIRA) can be used to verify claims of timber harvest location by matching levels of naturally occurring stable isotopes within wood tissue, to location-specific SIR predicted from reference data ('isoscaples'). The primary challenge in developing reliable isoscaples is the need to accurately predict stable isotopes in areas where no physical reference samples are available. Existing attempts to predict isoscaples from reference data have been hampered by the use of simple and ad-hoc statistical models, limiting the precision of estimated isoscaples and the confidence in derived estimates of geographical origin.

2: We present a new SIRA data analysis pipeline, designed to infer timber harvest location. We use Gaussian Processes to robustly estimate isoscaples from reference wood samples, which are then combined with species distribution range data to compute, for every pixel in the study area, the probability of it being the origin of the sample. Finally we present a methodology to determine priority locations to obtain new reference samples in future field expeditions.

3: We demonstrate our approach on a data set of  $n = 87$  wood samples from seven oak species in the USA as proof of concept. Our method is able to determine the harvest location up to 520-870 km, depending on the model parameterisation. Incorporating species distribution information improves accuracy by up to 36%. The new sampling locations proposed by our method decrease the variance of resultant isoscaples by up to 86% more than sampling the same number of locations at random.

4: The method we present here combines the prediction of isoscaples with derivation of geographical origin estimates. It advances the toolset available to authorities addressing illegal trade in forest products and enforcing anti-deforestation legislation. Importantly, reference data can be added as available, allowing for the expansion of reference collections and increasing prediction accuracy.

56 **Keywords:** SIRA, origin traceability, timber provenance, illegal logging, isoscaples,  
57 Gaussian Processes

## 58 1 Introduction

59

60

61

62

63

64

One million species now face extinction, and the unsustainable exploitation of natural resources is the second largest driver of terrestrial biodiversity loss, next only to land use changes [19]. To prevent our societies triggering a new wave of mass extinctions [3], nearly 200 nations have recently agreed on a new set of targets and goals under the Kunming-Montreal Global Biodiversity Framework. Under this international agreement, human-driven species extinctions must halt by 2030, in

65 order to allow an appropriate level of natural recovery by 2050. In particular, Target  
66 5 of the agreement has the objective to "ensure that the use, harvesting and trade  
67 of wild species is sustainable, safe and legal, preventing overexploitation".

68 Meeting such an ambitious but necessary target will require overcoming a key  
69 element of unsustainable use of natural resources: the illegal harvest of threatened  
70 trees. To combat illegal logging and associated trade in illegally harvested tim-  
71 ber, various frameworks have been put into place, such as the Convention on the  
72 International Trade in Endangered Species (CITES) and security linked sanctions.  
73 At national or regional levels, additional legislation include the US Lacey Act, the  
74 UK Timber Regulation, the EU Deforestation Regulation and the Australian Illegal  
75 Logging Prohibition Act. Despite the comprehensive legal framework already in  
76 place, and the international commitments under current adoption, it is clear that  
77 none of these can effectively work without enforcement. This is where technological  
78 solutions and methodological developments become key.

## 79 **1.1 Stable Isotope Ratio Analysis for timber provenance**

80 Current legislation increasingly requires accurate, cost-effective, and high-throughput  
81 tools that can verify the specific species and origin of products in global trade [20].  
82 Scientific testing technologies have become relatively well established and are al-  
83 ready supporting companies and enforcement authorities to scrutinize traceability  
84 systems. These technologies measure the chemical, anatomical and genetic features  
85 of plants which, when compared against a robust physical reference collection, can  
86 be used to (in-)validate declared species and origin claims and support enforcement  
87 officials in their efforts to detect, for example, illegal or conflict timber and fraud in  
88 supply chains.

89 One of the most promising and widely used scientific technologies is stable iso-  
90 tope ratio analysis (SIRA). The ratios of several elemental stable isotopes within  
91 natural products vary across space and can assist in verifying the geographic origin  
92 of products. Most of the abundant elements in organic compounds (Hydrogen,  
93 Oxygen, Carbon, Sulfur, Nitrogen) have naturally-occurring stable isotopes that do  
94 not undergo radioactive decay, and can be readily detected by mass spectrometry  
95 [5]. The composition of stable isotopes incorporated into the tissues of a plant is  
96 determined by the soil, climate, metabolic fractionation and other biotic and abi-  
97 otic conditions characteristic of the species and its habitat [41, 9, 26]. Stable  
98 isotope ratios can be used to discriminate between geographic areas as the varia-  
99 tion in these stable isotope ratios depend on natural variations of the underlying  
100 mechanisms (for example environmental drivers [47]). The types of products that  
101 can be analysed by SIRA to determine risk for being illegally harvested include  
102 natural resources-such as forest products [48, 5], agricultural products [10, 39],  
103 wildlife [8, 32], fish/seafood [16, 42, 34], ivory [46, 54], precious metals [31] and  
104 illicit drug trafficking, including natural and synthetic opioids [35, 11], and other  
105 forensic uses to identify counterfeit and pirated goods trafficking [12].

## 106 1.2 Modelling approach

107 Current modelling practices for the use of SIRA to verify harvest location of both  
108 legally and illegally harvested products could be improved. The use of SIRA is  
109 currently limited by the relative simplicity of models used, as well as by the limited  
110 number of reference samples used as input data for such models. The practical  
111 nature of reference sampling campaigns is that they can be costly and budgetary  
112 needs are often underestimated, with sampling locations often taking into account  
113 relative ease of sampling rather than areas that yield a gain in model prediction  
114 accuracy [40]. There has been considerable development of isoscapes ("isotope  
115 landscapes"), given that stable isotopic variation is a continuous spatial variable in  
116 nature [50]. These isoscapes are geospatial maps that show the isotopic variation  
117 of the material of interest [50]. While the potential of isoscapes for determining  
118 product origins has long been recognized, few rigorous methods exist to achieve  
119 this task. The existing methods use simple prediction strategies such as linear  
120 regression [48, 49] and clustering [32], which do not fully leverage the information  
121 contained in isotope ratio data.

122 Gaussian Process (GP) regression, also known as kriging in geostatistics litera-  
123 ture, is a class of flexible regression models which use the values in sampled points  
124 to estimate the values in surrounding, unsampled points [36, 18, 51]. A key advan-  
125 tage of GP regression is that it can quantify the uncertainty of its own predictions  
126 based on the inferred spatial covariance structure of the samples. Quantifying the  
127 uncertainty of predictions is viewed as increasingly important in safety-critical [28]  
128 and forensic [44, 45] machine learning applications. GP regression also facilitates  
129 co-kriging: inferring the values of a sparsely sampled variable of interest through  
130 variables that are highly correlated with it but more densely sampled [2, 30]. In the  
131 context of plant origin estimation, co-kriging translates to inferring stable isotope  
132 ratios from atmospheric drivers (such as precipitation, temperature and water vapor  
133 pressure) known to influence the stable isotope signal in wood [26, 41]. This then  
134 provides a powerful tool for predicting the isotopic composition in areas that have  
135 not yet been sampled (examples are given in [23, 48]).

136 Previous work on isoscapes used GP regression primarily as a spatial interpo-  
137 lation technique without a probabilistic interpretation [23, 32, 30]. A more recent  
138 method uses GP variance estimates from precipitation isoscapes for origin deter-  
139 mination in animals [37]. In this work, we develop GP-based probabilistic machine  
140 learning models to infer the origin of timber samples by directly modelling timber  
141 isoscapes. We present a new data analysis pipeline that incorporates timber SIRA  
142 data, atmospheric predictors and species distribution data. We find that probabilis-  
143 tic modeling greatly enhances the utility of SIRA in estimating the geographical  
144 origin of timber and helps guide future sample collection by identifying sampling  
145 locations that will minimize prediction uncertainty. The presented framework can  
146 then be applied to trace back timber of endangered species, by assisting in deter-  
147 mining where to collect samples and by using SIRA datasets being collected by the  
148 World Forest ID [22] initiative.

## 149 2 Materials and Methods

### 150 2.1 Data sets

151 We use data from 87 trees of the genus *Quercus*, sampled across the contiguous  
152 United States, as described in Watkinson et al. (2020) [48]. Stable isotope ratio  
153 measurements were done following the protocol described in Boner et al. (2007) [5]  
154 and Watkinson et al. (2020) [48]. Each entry contained stable isotope ratio mea-  
155 surements of oxygen  $\delta^{18}\text{O}$  (ratio between  $^{18}\text{O}$  and  $^{16}\text{O}$ ), hydrogen  $\delta^2\text{H}$  (ratio be-  
156 tween  $^2\text{H}$  and  $^1\text{H}$ ), carbon  $\delta^{13}\text{C}$  (ratio between  $^{13}\text{C}$  and  $^{12}\text{C}$ ) and sulfur  $\delta^{34}\text{S}$  (ratio  
157 between  $^{34}\text{S}$  and  $^{32}\text{S}$ ) as well as the GPS coordinates of the sampled tree. Stable  
158 isotope ratios are largely driven by environmental conditions such as precipitation,  
159 temperature, humidity and so on. Thus, it is natural to use publicly available data  
160 on those factors to aid the inference of isoscapes. We used the following atmo-  
161 spheric data:  $\delta^2\text{H}$  and  $\delta^{18}\text{O}$  isotopic composition of precipitation [7], water vapour  
162 [6] (found to be associated with  $\delta^{13}\text{C}$  by Watkinson et al. [48]), reflected shortwave  
163 radiation [1] and precipitation (multi-satellite) [27], both of which were found to be  
164 associated with  $\delta^{34}\text{S}$  by Watkinson et al. [48].

165 To inform the priors (probability distributions representing the prior belief on  
166 possible tree locations) of the models we develop, we used species inventory data  
167 across the natural range of each species within the United States [52], downloaded  
168 from: <https://www.fs.usda.gov/rds/archive/Catalog/RDS-2013-0013> on  
169 09/12/2022. This data is available as species-specific raster layers of tree abundance  
170 at 250m resolution. We then used the function `project()` of the R package *terra* [25]  
171 to bilinearly aggregate abundance data so that it matched the spatial resolution of  
172 other spatial data in the pipeline.

### 173 2.2 Model architecture

174 Figure 5 presents the pipeline overview of the data sets, components comprising our  
175 model and output. We use a rectangular grid to represent the study area. Grid points  
176 are placed every 0.125 degree latitude ( $\approx 14$  km) and every 0.06 degree longitude  
177 ( $\approx 4.3$ -6.0 km), which allows us to approximate spatial probability distributions  
178 with high accuracy. For every isotope ratio (IR), we fit a GP regression model to  
179 the training data to obtain the posterior mean and variance of the isotope ratio for  
180 every point of the grid - see Sections 2.3 and 2.4 for more details. The input to  
181 the GP consists of the coordinates and/or the climate variable values at the grid  
182 point. For a combination of stable isotope ratios  $(y_O, y_H, y_C, y_S)$  of observed stable  
183 isotope ratios, we compute the likelihood of observing it at every point in the grid  
184 using the four GP regression models estimated in the previous step. This likelihood  
185 is the product of likelihoods for each isotope ratio as we assume independence  
186 between isotopes. Given the prior and the likelihood, we perform Bayesian inference  
187 by computing the posterior probability of each grid point being the origin of the  
188 sample by applying Bayes' Theorem. For ease of interpretation, the output is a map  
189 with highest-posterior density (HPD) regions indicated for several probability levels  
190 (15%, 30%, 50%, 75%, 90%, 95%).

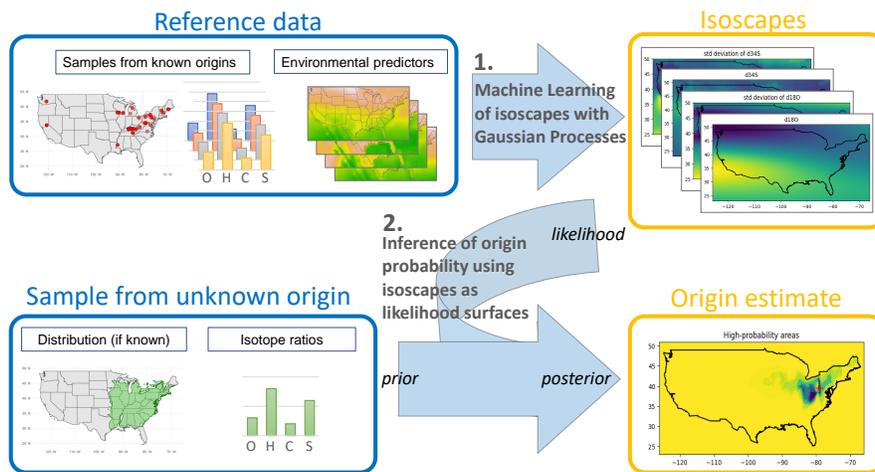


Figure 1: Model workflow. We use a training set of isotope ratios from trees collected at known locations and atmospheric data layers ("Reference data"). We fit a Gaussian process regression model to infer isoscapes and associated variance estimates, and compute the likelihood of observing the IR value for each element across the study area. To estimate the source of material with uncertain provenance ("Samples from unknown origin"), the isoscapes are then combined with prior information on the geographical distribution of the species, to yield a probability distribution of origin for the sample. We visualize predicted probability maps by plotting highest-posterior density regions for several probability levels (15%, 30%, 50%, 75%, 90% and 95%, dark blue to light green).

191 **2.3 Gaussian Process regression**

192 In the following, we give a brief overview of GPs. For a more thorough explanation,  
 193 see [51].

194 GPs provide a flexible framework for regression, which enables the modeller to  
 195 quantify the uncertainty of specific inferences. A GP is a random process such that  
 196 all of its marginals are jointly normally distributed (Gaussian). Let  $\mathbf{x} = [x_{lon}, x_{lat}]$   
 197 be the GPS coordinates of a sample. For any set of positions  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ ,  
 198 the responses  $y_1, y_2 \dots, y_n$  at those positions are assumed to be jointly normally  
 199 distributed

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} m(\mathbf{x}_1) \\ m(\mathbf{x}_2) \\ \vdots \\ m(\mathbf{x}_n) \end{bmatrix}, \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & k(\mathbf{x}_n, \mathbf{x}_2) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix} + \sigma^2 \mathbf{I} \right) \quad (1)$$

200 where:

- 201 1.  $m(\mathbf{x})$  is the *mean function* describing the *a priori* expected value of the re-  
 202 sponse  $y$  at point  $\mathbf{x}$ . Usually, this is a parameterized function whose pa-  
 203 rameters are estimated by fitting the model to the training data. The most  
 204 common choice of mean function is the constant mean  $m_c(\mathbf{x}) = \theta_m$  for all  
 205  $\mathbf{x}$ , which we also use in this work.
- 206 2.  $k(\mathbf{x}_1, \mathbf{x}_2)$  is the *covariance function* describing the *a priori* covariance be-  
 207 tween responses at points  $\mathbf{x}_1$  and  $\mathbf{x}_2$ . This is also a parameterized func-  
 208 tion. Popular choices of  $k$  are the *squared exponential*  $k_{se}(\mathbf{x}_1, \mathbf{x}_2) =$   
 209  $A \exp(-|\mathbf{x}_1 - \mathbf{x}_2|^2 / \rho^2)$ , or the Matern function [51], which both reflect the  
 210 common assumption that similar predictor values will lead to similar response  
 211 values. In this work, we use the Matern function with separate scaling pa-  
 212 rameters for latitude and longitude to model spatial covariance.
- 213 3.  $\sigma^2$  is the intrinsic noise parameter
- 214 4.  $\mathbf{I}$  is the  $n \times n$  identity matrix.

215 We will write  $\mathbf{y}$ ,  $\mathbf{m}$  and  $K$  to denote the responses, means and the covariance matrix  
 216 of the training data, respectively, so that we can write Eq. 1 as  $\mathbf{y} \sim \mathcal{N}(\mathbf{m}, K + \sigma^2 \mathbf{I})$ .  
 217 The choice of mean and covariance functions reflects prior knowledge and modelling  
 218 assumptions about the regression problem. The function parameters as well as  
 219 the noise parameter  $\sigma$  are estimated by maximizing the likelihood of the training  
 220 data. We use GPyTorch [21] to efficiently find the maximum likelihood parameter  
 221 estimates.

222 After parameter estimation, the GP regression model can be used to predict  
 223 responses at previously unseen data points. Let  $S$  be the locations and responses  
 224 comprising the training data set. Since the responses at training and test points  
 225 are assumed to be jointly Gaussian, the conditional distribution of the response at  
 226 a test point  $\mathbf{x}^*$  given the training data is also Gaussian with mean

$$\mu(\mathbf{x}^*|S) = m(\mathbf{x}^*) + \mathbf{k}^*(K + \sigma^2\mathbf{I})^{-1}(\mathbf{y} - \mathbf{m}) \quad (2)$$

227 where  $\mathbf{k}^* = [k(\mathbf{x}^*, \mathbf{x}_1), k(\mathbf{x}^*, \mathbf{x}_2), \dots, k(\mathbf{x}^*, \mathbf{x}_n)]$  is the vector of a priori covariances  
 228 between responses at  $\mathbf{x}^*$  and training data points. The posterior variance of  $y^*$  is  
 229 given by

$$\sigma^2(\mathbf{x}^*|S) = k(\mathbf{x}^*, \mathbf{x}^*) + \sigma^2 - \mathbf{k}^*(K + \sigma^2\mathbf{I})^{-1}\mathbf{k}^{*\top} \quad (3)$$

230 - see [51] for a derivation.

231 For a specific response value  $y^+$ , its likelihood at  $\mathbf{x}^*$  is just the Gaussian prob-  
 232 ability density with mean  $\mu$  and variance  $\sigma^2$  found by applying Equations 2 and  
 233 3

$$p(y^* = y^+|\mathbf{x}^*, S) = \frac{1}{\sqrt{2\pi\sigma^2(\mathbf{x}^*|S)}} \exp\left(\frac{-(y^+ - \mu(\mathbf{x}^*|S))^2}{2\sigma^2(\mathbf{x}^*|S)}\right) \quad (4)$$

## 234 2.4 Incorporating atmospheric data

235 For any  $\mathbf{x}$ , let  $\mathbf{u}_i(\mathbf{x})$  denote the 12-entry vector of monthly values of atmospheric  
 236 variable  $i$  at location  $\mathbf{x}$ . We use a linear covariance function to model the covariance  
 237 component corresponding to the variation in atmospheric variable  $i$

$$k_i(\mathbf{x}_1, \mathbf{x}_2) = \theta_i[\mathbf{u}_i(\mathbf{x}_1)]^\top \mathbf{u}_i(\mathbf{x}_2) \quad (5)$$

238 where  $\theta_i$  is a parameter to be estimated during training. The linear covariance  
 239 function models a linear relationship between the atmospheric variable and the  
 240 response and is mathematically equivalent to Bayesian linear regression [51].

241 The overall covariance function is the sum of the spatial term and the linear  
 242 terms

$$k(\mathbf{x}_1, \mathbf{x}_2) = k_{spatial}(\mathbf{x}_1, \mathbf{x}_2) + \sum_{i \in V} \theta_i[\mathbf{u}_i(\mathbf{x}_1)]^\top \mathbf{u}_i(\mathbf{x}_2) \quad (6)$$

243 where  $V$  is the set of atmospheric variables impacting the considered isotope ratio.

## 244 2.5 Bayesian inference of sample origin

245 Given a prior distribution  $p(\mathbf{x})$  over possible sample origins and a GP regression  
 246 model for every isotope, we perform Bayesian inference of sample origin. For a  
 247 sample  $\bar{y} = (y_O, y_H, y_C, y_S)$  of observed values, the Bayes' theorem gives the  
 248 posterior distribution of possible origins:

$$p(\mathbf{x}|\bar{y}, S) = \frac{p(\mathbf{x}) \prod_{i \in \{O, H, C, S\}} p_i(y_i|\mathbf{x}, S)}{\int_{\mathbf{x} \in A} p(\mathbf{x}) \prod_{i \in \{O, H, C, S\}} p_i(y_i|\mathbf{x}, S) d\mathbf{x}} \quad (7)$$

249 where the probabilities  $p_i$  are computed from the GP models for the respective  
 250 isotopes using Equation 4 and  $A$  is the study area. The integral in the denominator  
 251 is computed by averaging the probabilities over the spatial grid.

## 252 **2.6 Incorporating species distributions**

253 We use the spatial density maps developed by Wilson *et al.* [52] to design two  
254 prior distributions for sample origin that account for the spatial distribution of oak  
255 species. The first, which we call the *density prior*, holds that the probability of a  
256 sample originating at a grid cell is proportional to the basal area (average amount  
257 of area occupied by tree stems per unit of space) recorded at the grid cell. The  
258 second, which we call the *range prior*, assigns equal probability to every grid cell  
259 where above-zero basal area has been recorded. In addition, both priors allow for a  
260 small probability that a sample might occur outside its observed range - we set that  
261 probability to 0.01 and diffuse it uniformly over all grid points within the contiguous  
262 United States where the species does not occur according to Wilson *et al.*.

## 263 **2.7 Performance evaluation**

264 We perform 5-fold cross-validation on the data set and report the average values of  
265 all performance metrics over all data points. Samples with incomplete or ambigu-  
266 ous species information and samples collected in botanical gardens outside of their  
267 species' native range are excluded from the test sets, but not from the training sets,  
268 resulting in a total of n=74 test samples across the 5 folds.

269 Rigorously evaluating the performance of our model is a non-trivial task as it  
270 produces a distribution over possible locations, rather than a single location. For  
271 this reason, we have defined several metrics to investigate different aspects of our  
272 predictions.

### 273 **2.7.1 Predictive log-likelihood and log-posterior**

274 We report the log-likelihood (Eq. 4) and the log-posterior (Eq. 7) of observing the  
275 sample  $\bar{y}$  at its true origin  $\mathbf{x}_t$ . Both of those measure how well the model fits the  
276 test data.

### 277 **2.7.2 Mean posterior distance to true location (MPD)**

To investigate how distant our predicted locations are from the truth, we report the  
expected distance between the true location and a location sampled randomly from  
the posterior distribution returned by our model

$$MPD = \int_{\mathbf{x} \in A} d(\mathbf{x}_t, \mathbf{x}) p(\mathbf{x} | \bar{y}, S) d\mathbf{x}$$

278 where  $d()$  is the great circle distance between the two points. A perfect prediction  
279 would have the distance of 0. In practice, isoscapes often predict similar isotope  
280 ratio values at distant locations, so even a statistically efficient method might yield  
281 a high MPD value. This metric will favour predictions concentrated around the  
282 true location over equally dispersed predictions concentrated elsewhere. It will also  
283 favour less dispersed predictions generally.

284 **2.7.3 Area scored higher than the true location (ASH)**

The behaviour of MPD is influenced by the shape of the posterior distribution, which favours unimodal over multimodal shapes. We report the total surface area corresponding to the points that the model considers more plausible than the true origin of the sample.

$$ASH = \int_{\mathbf{x} \in A} I[p(\mathbf{x}|\bar{y}, S) > p(\mathbf{x}_t|\bar{y}, S)] d\mathbf{x}$$

285 where  $I(\cdot)$  is the indicator function that yields 1 when the statement is true and 0  
 286 otherwise. In contrast to MPD, this metric is likely to give a low value to a posterior  
 287 distribution that is concentrated in several small areas as long as one of those areas  
 288 contains the true location.

289 **2.8 Guiding future sampling efforts**

290 Field sample collections are time-consuming and expensive. By having an idea  
 291 which sample points need to be collected for increased origin estimation accuracy,  
 292 we can optimize future field collections. The isoscape variance estimates provided  
 293 by GPs can be used to guide future sampling efforts, which in turn will maximize  
 294 the performance of the model. This approach is known as *active learning* in the  
 295 machine learning literature. Here, we propose a strategy to minimize the error of  
 296 our isoscape estimates by carefully choosing future sampling locations.

297 Early attempts at efficient active learning in GPs involved collecting samples  
 298 at points with highest response variance or, equivalently, picking a set of points  
 299 that maximizes the entropy of responses [15]. Unfortunately, this approach tends  
 300 to recommend collecting samples on the boundaries of the study area, which is  
 301 wasteful as the newly collected samples improve isoscapes in a smaller fraction of  
 302 the study area than if they were placed away from the boundary. This motivated  
 303 researchers to propose several criteria for optimizing sampling [24, 38]. Here, we  
 304 adopt an approach similar to that of Guestrin *et al.* [24] with a few modifications  
 305 designed to address the large size of our spatial grid, which renders their original  
 306 method computationally intractable for our data set.

307 We seek to maximize the *average* reduction in predictive variance across our  
 308 study area that can be achieved by adding a sample to the training set. Let  $S$   
 309 be the set of sampled locations and  $G$  be the set of grid points. We define the  
 310 information gain (IG) associated with adding a new point  $\mathbf{x}^*$  to the training data  
 311 set as

$$IG(\mathbf{x}^*) = \sum_{\mathbf{x} \in G} [(\log(\sigma^2(\mathbf{x}|S)) - \log(\sigma^2(\mathbf{x}|S \cup \{\mathbf{x}^*\})))] \quad (8)$$

312 where the predictive variances are computed using Equation 3. The algorithm then  
 313 picks the point in the grid that yields the highest IG. Importantly, the predictive  
 314 variances depend only on the sampling locations and not on the sampled values,  
 315 so it is possible to sequentially propose multiple sampling points before collecting  
 316 the samples. Our method sequentially proposes sample collection points until a  
 317 user-specified number of samples is reached. We assume that samples can only

318 be collected in locations where at least one of the species is present. Thus, grid  
319 points that lie outside every species range are excluded from the procedure. Our  
320 active learning strategy requires repeatedly computing a large number of predic-  
321 tive variances for varying training sets. To reduce computation time, we randomly  
322 downsample our grid to 15000 points before running the analysis. In addition, we  
323 assume that the reduction in variance is negligible for grid points situated more than  
324 15 degrees away from  $x^*$  in longitude or more than 7.5 degrees in latitude.

## 325 **3 Results**

### 326 **3.1 Model accuracy and comparison**

327 Table 1 shows performance metrics for all the models on the test data set. In  
328 all settings, the plausible origin areas identified by our models consisted of points  
329 within an average distance of 520-870 kilometers from the true locations. Even  
330 with a relatively small training data set of 69 training samples, our model is able  
331 to exclude the vast majority of the study area as a possible source of the sample  
332 under consideration.

333 Incorporating species distribution information improves prediction performance  
334 for every model and every metric examined except the log-likelihood, which does  
335 not depend on the prior. Informative priors improve MPD by 16% to 35% and ASH  
336 by 15% to 57% with most improvement for the pure spatial model and least for  
337 the spatial+atmospheric model. The more informative density prior gives better  
338 accuracy than the range prior according to all the metrics. Predicted probability  
339 maps for a few test points are shown in Fig. 2 and 3.

340 The spatial-only GP model gives the closest predictions to the true location,  
341 except when a flat prior is used. In general, the spatial-only and the combined  
342 spatial+atmospheric model give similar results on all metrics and they both outper-  
343 form the atmospheric-only model in almost all settings. Somewhat surprisingly, the  
344 combined model does not outperform the spatial-only model. This might be due to  
345 the relatively small data set size.

346 The predictions of atmospheric GP models appear qualitatively different from  
347 those from the purely spatial GP. Atmospheric model predictions often emphasize  
348 geographical areas with distinct climate patterns, such as the Appalachia or the Gulf  
349 Coast. Unsurprisingly, the purely spatial GP identifies areas that are more spatially  
350 cohesive but do not share any obvious physical features.

### 351 **3.2 Guiding future sampling efforts**

352 We investigated the performance of our active learning strategy on the US oak  
353 data set. For the spatial-only model, we let our method propose  $n_s = 10$  new  
354 sampling locations to add to the training data set in the first cross-validation fold and  
355 computed the predictive variances before and after including the proposed locations.

356 The resulting isoscape standard deviation maps are shown in Figure 4. Our active  
357 learning strategy proposes sampling locations in currently undersampled regions with

model	prior	log L	MPD (km)	log-posterior	ASH (km <sup>2</sup> )
Spatial-only	flat	<b>-6.964</b>	809	-9.582	470000
Spatial-only	range	<b>-6.964</b>	600	-9.537	327000
Spatial-only	density	<b>-6.964</b>	<b>520</b>	-9.059	<b>203000</b>
Atmospheric-only	flat	-7.362	870	-9.972	576000
Atmospheric-only	range	-7.362	606	-9.797	450000
Atmospheric-only	density	-7.362	567	-9.428	311000
Atmospheric+Spatial	flat	-7.149	794	-9.518	382000
Atmospheric+Spatial	range	-7.149	627	-9.431	315000
Atmospheric+Spatial	density	-7.149	536	<b>-8.978</b>	213000

Table 1: Mean test set performance for all the models used in the study. Best values across all models are shown in bold. The Spatial-only GP combined with the density prior gives the highest predictive log-likelihood and log-posterior and the lowest MPD and ASH values for all priors used. The Spatial-only model outperforms the other models when range or density priors are used, while the Atmospheric+Spatial model performs best in terms of MPD and ASH when flat priors are used. The inclusion of species distribution information decreases MPD and ASH values for all models used.

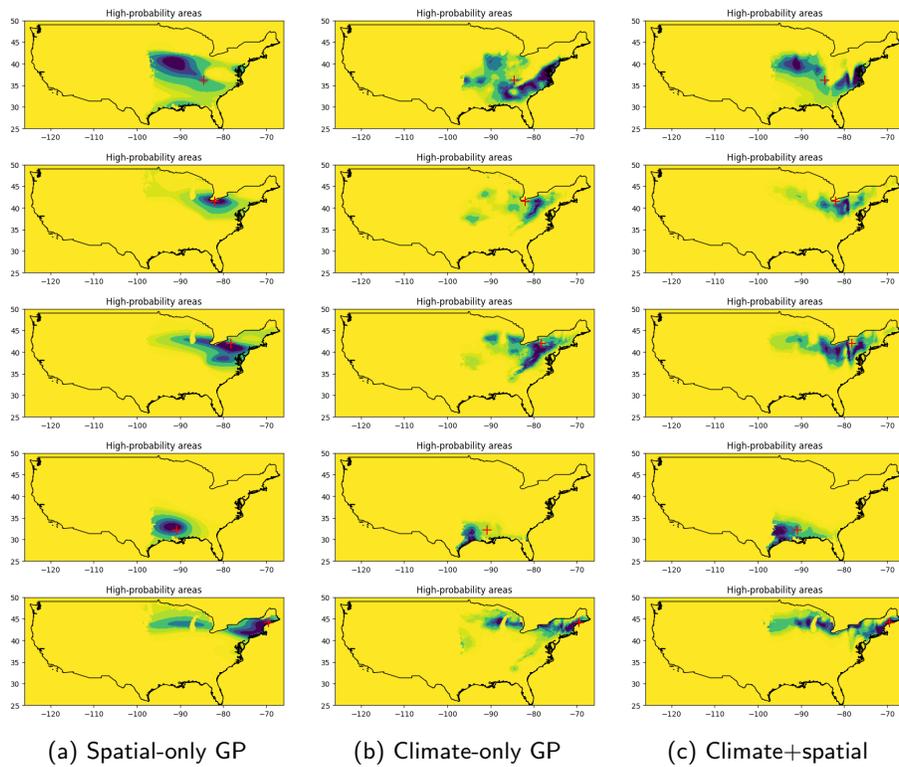


Figure 2: Spatial predictions from the three models for 5 points from the test set using the range prior. Darker shades denote areas with higher probability mass and the red cross indicates the actual location of the tree.

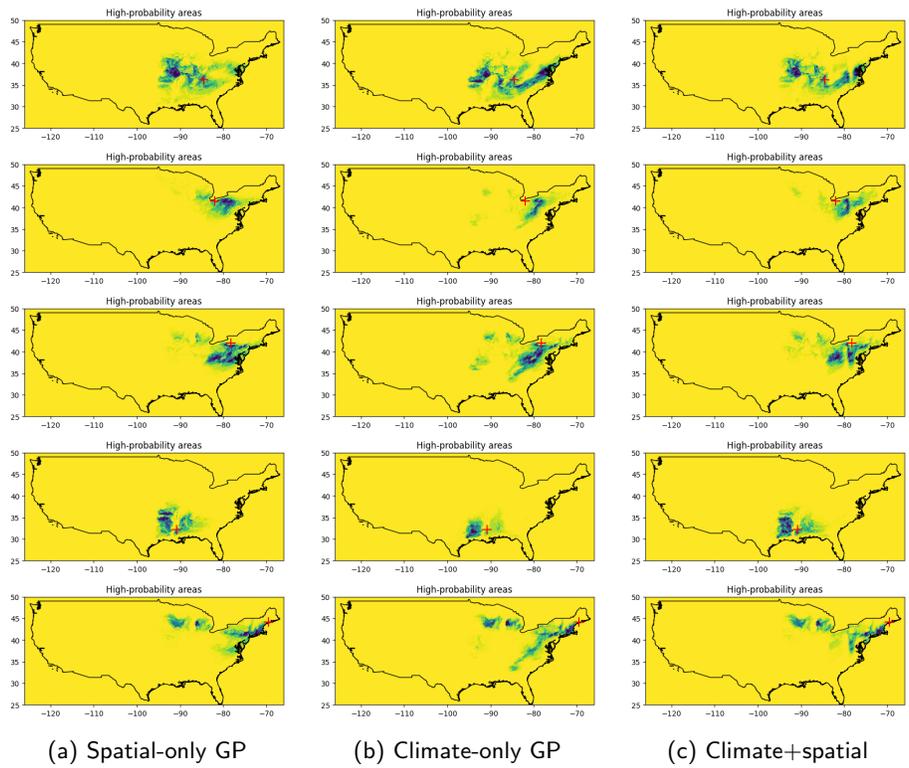


Figure 3: Spatial predictions from the three models for 5 points from the test set using the density prior. Darker shades denote areas with higher probability mass and the red cross indicates the actual location of the tree.

358 high predictive variance and sampling in those areas results in a visible improvement.  
 359 The highest decrease in predictive variance was observed for  $\delta^2\text{H}$  while the lowest  
 360 decrease was observed for  $\delta^{18}\text{C}$ . Most of the chosen locations are close to, but not  
 361 at the boundary of the allowed sampling area.

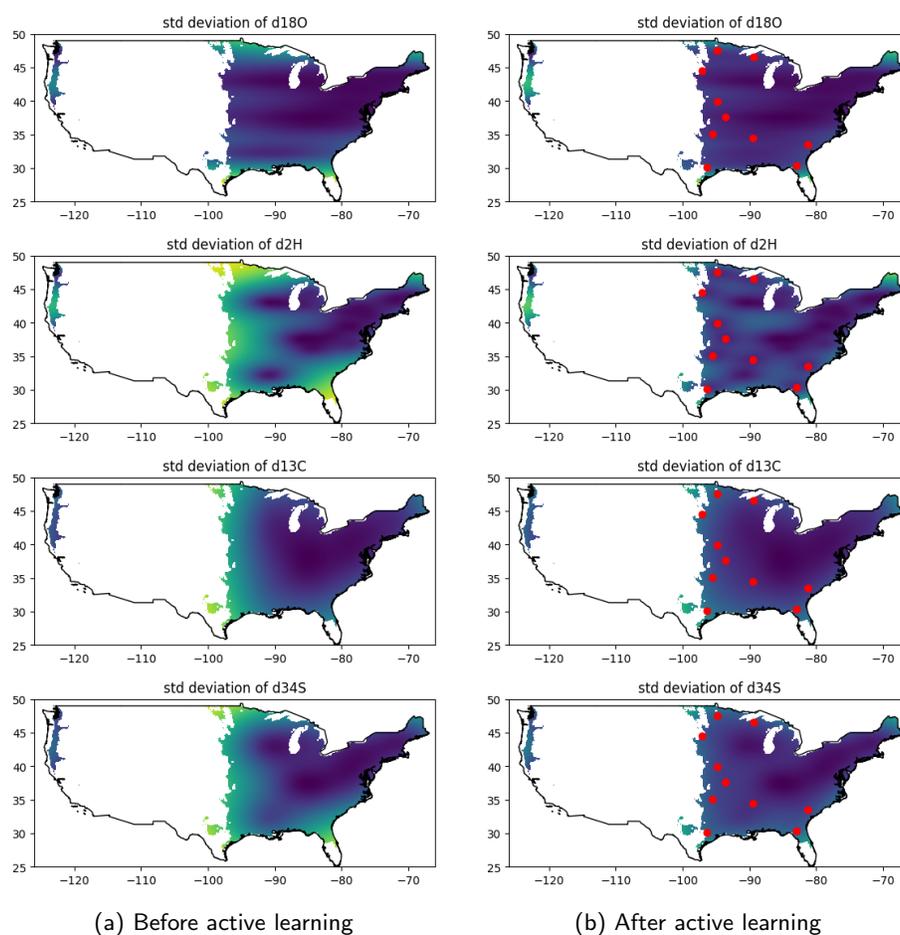


Figure 4: Maps showing predictive standard deviations for the four isotopes before and after adding  $n_s = 10$  sample locations proposed by our active learning method for the spatial-only model. Standard deviations are only shown within the allowed sampling area, which is the union of ranges for the species in our data set. The red dots show the proposed locations. Our method proposes locations in areas with high predictive variance, particularly for  $\delta^2\text{H}$  and  $\delta^{34}\text{S}$ . Adding the proposed locations leads to a marked reduction of variance in the neighboring areas.

362 To investigate the efficiency of our active learning procedure, we compared  
 363 isoscape variances resulting from active learning with those resulting from adding

364 the same number of points sampled randomly from the allowed sampling area. We  
365 generated  $n_r = 100$  such variance maps and compared the average variance (across  
366 the allowed sampling area) of those maps with the maps in Fig. 4. Fig.S1 shows the  
367 average predictive variances as a function of the number of points added for both  
368 random and active learning sampling strategies. We see that our active learning  
369 strategy results in substantially faster decrease in predictive variances. After adding  
370 10 samples, the reduction in variance achieved by our active learning method is  
371 between 64% ( $\delta^{13}C$ ) and 86% ( $\delta^{18}O$ ) greater than the average reduction achieved  
372 by the same number of random samples.

## 373 4 Discussion

### 374 4.1 Timber origin estimation

375 To halt illegal logging, to enforce timber regulations and to protect biodiversity in  
376 forested landscapes, we need to be able to accurately estimate the timber harvest lo-  
377 cation. Although several examples exist of applying SIRA for timber origin questions  
378 [23, 48, 29], these approaches do not take full advantage of (1) atmospheric and  
379 species distribution datasets available or (2) state-of-the-art probabilistic machine  
380 learning models. In this work we present a new computational pipeline which aims at  
381 taking advantage of both. The accuracy of our models depends on the specific mod-  
382 elling approach being used and the data sets incorporated. Using prior information  
383 about species distribution results in a considerable increase in accuracy regardless  
384 of which model is used by all metrics considered. The impact of adding species  
385 distribution data appears to be greater for the spatial-only model than models that  
386 use atmospheric information. This could be due to climate patterns influencing both  
387 species distributions (habitat suitability) and the values of the atmospheric variables  
388 that we incorporated in our models, which renders species distribution information  
389 more redundant once atmospheric variables have been included in the model.

390 Within timber tracing literature, our method bears the most resemblance to the  
391 work of Watkinson *et al.* [48], which uses linear regression to predict isoscapes based  
392 on atmospheric data. Their approach assumes a constant variance across the study  
393 area. In contrast, our method estimates the predictive variances based on the spatial  
394 covariance structure learned from the reference data, which enables us to translate  
395 differences in sampling density across regions into varying levels of confidence in  
396 isoscapes across space. Like Watkinson *et al.* [48], our method assumes a linear  
397 relationship between atmospheric predictors and isoscapes, but our GP formulation  
398 implicitly integrates over plausible values of regression parameters, which should lead  
399 to more robust predictions compared to standard linear regression. In addition, our  
400 approach makes use of species distribution data, which yields substantially improved  
401 predictions compared to uninformative priors. Finally, our approach enables us to  
402 propose locations for further sample collection that maximize the utility of the  
403 samples.

404 Estimating the spatial covariance structure has recently attracted attention in  
405 animal stable isotope studies. Ma *et al.* [37] recently proposed a method that uses

406 probabilistic precipitation isoscapes derived from a GP [14], which are then cali-  
407 brated to produce isoscapes for the species of interest. St. John Glew *et al.* [43]  
408 introduced a model combining spatial and environmental effects using a novel likeli-  
409 hood approximation for isoscape estimation, though the main focus of their work is  
410 isoscape modelling, not origin estimation. These approaches differ from ours in that  
411 1) they rely on Laplace approximations for isoscape estimation rather than exact  
412 likelihood maximization; 2) they use ordinary least-squares regression to account for  
413 atmospheric predictors, whereas our method uses a Bayesian approach via a linear  
414 covariance term; and 3) they do not aim to actively improve isoscapes through  
415 additional sampling. A common feature between these models and ours is using a  
416 grid to compute the posterior distribution of origins, which was first considered by  
417 Wunder [53].

418 Our current best performing model can estimate the origin of harvest location  
419 for *Quercus* species to 520 km across the (north-)east of the United States. Future  
420 field expeditions will lead to an improvement, especially if the identified priority  
421 locations are targeted (see 4.2). The presented model will be adapted to other  
422 use cases, with mainly a focus on tropical species on which the logging pressure is  
423 significant and which might be endangered.

## 424 **4.2 Guiding future collection efforts**

425 We expect that our models will be more accurate once more timber samples be-  
426 come available. The size of the current data set of wood samples available to this  
427 study ( $n=87$ ) is quite small relative to the area of contiguous United States, which  
428 inevitably results in large predictive variance in many areas. In addition to reducing  
429 uncertainty about undersampled areas, larger data sets (in the range of hundreds to  
430 thousands of samples collected from across the US) should also enable researchers  
431 to use more complex GP models, including models with heterogeneous noise [4], or  
432 deep GP models where the covariance function is modelled by a neural network [17].

433 Under the World Forest ID Programme [22], tens of thousands of tree samples  
434 are being collected globally, and are being analysed by different techniques, including  
435 SIRA. Our active learning approach can be used to inform future sample collection  
436 efforts and increase model accuracy that can be achieved within a fixed sampling  
437 budget. This will be especially important in tropical regions, where reaching sam-  
438 pling sites can be difficult, time intensive and expensive. A good sampling design  
439 can substantially improve model performance [13], and our method can be used to  
440 adapt sampling efforts as more data is analysed. Our current approach focuses on  
441 minimizing predictive variances without considering the impact of newly sampled  
442 points on model parameters. Extending our approach to *non-myopic* sampling [33],  
443 which considers the impact on model parameters, would constitute an interesting  
444 future research direction. Another avenue for improving our approach would be  
445 to augment our IG criterion to reflect the varying cost of collecting samples as a  
446 function of the time and financial cost of reaching the desired sampling location.

## 447 **5 Conclusion**

448 The accurate estimation of geographic origin of globally traded wood products  
449 is a critical step in combating illegal logging and associated trade, by supporting  
450 authorities' ability to verify claims made by traders at any supply chain node. In this  
451 work we presented a novel analytical pipeline that brings together and incorporates  
452 multiple data types and algorithms. This methodology is able to accurately predict  
453 timber product origin and can be used to optimize future field sampling to further  
454 increase accuracy and precision. We hope that this work will inspire more efforts  
455 to expand reference collections of wood samples, such as under the auspices of the  
456 World Forest ID Programme (<https://worldforestid.org/>), and that governments  
457 and companies will more routinely use the technological tools at their disposal to  
458 have more oversight over their supply chains and promote a more sustainable use  
459 of natural resources.

## 460 **6 Conflict of interest statement**

461 The authors declare that they have no conflicts of interest.

## 462 **7 Data availability**

463 The data and code used in this study will be made publicly available within a year  
464 of publication. Earlier access may be provided by request.

## 465 **8 Author contributions**

466 Jakub Truszkowski, Victor Deklerck and Alexandre Antonelli jointly conceived the  
467 project. Jakub Truszkowski designed the methodology, implemented the algorithms,  
468 performed most of the analyses and wrote parts of the manuscript. Roi Maor  
469 selected and pre-processed some of the data sets, designed the main figure and wrote  
470 parts of the manuscript. Raquib Bin Yousuf, Subhodip Biswas, Naren Ramakrishnan  
471 and John Simeone wrote some of the software code and performed initial data  
472 explorations. Scot McQueen selected the data sets for species distribution models.  
473 Caspar Chater, Peter Gasson, Marigold Norman and Jade Saunders wrote parts  
474 of the manuscript. Victor Deklerck directed the project and wrote parts of the  
475 manuscript. All authors gave final approval for submission.

## 476 **9 Acknowledgements**

477 Jakub Truszkowski and Alexandre Antonelli are funded by the Swedish Research  
478 Council (grant number 2019-05191). Victor Deklerck is funded under the World  
479 Forest ID Timber at Kew Grant provided by the Department of Environment, Food  
480 & Rural Affairs (DEFRA), International Climate Finance (ICF) R&D Programme,

481 UK (project 29084). Caspar Chater and Roi Maor are funded under the World  
482 Forest ID FRC at Kew grant provided by DEFRA, ICF R&D Programme, UK.  
483 Alexandre Antonelli also acknowledges financial support from the Swedish Foun-  
484 dation for Strategic Environmental Research MISTRA (Project BioPath) and the  
485 Royal Botanic Gardens, Kew. The authors want to thank the US Forest Service -  
486 International Programs and FSC-US for the initial collection of the dataset. The  
487 findings and conclusions in the article are those of the authors.

## 488 References

- 489 [1] URL: [https://neo.gsfc.nasa.gov/view.php?datasetId=CERES\\_](https://neo.gsfc.nasa.gov/view.php?datasetId=CERES_)  
490 SWFLUX\_M.
- 491 [2] Sajal Kumar Adhikary, Nitin Muttill, and Abdullah Gokhan Yilmaz. Cokriging  
492 for enhanced spatial interpolation of rainfall in two australian catchments.  
493 *Hydrological processes*, 31(12):2143–2161, 2017.
- 494 [3] Anthony D Barnosky, Nicholas Matzke, Susumu Tomiya, Guinevere OU  
495 Wogan, Brian Swartz, Tiago B Quental, Charles Marshall, Jenny L McGuire,  
496 Emily L Lindsey, Kaitlin C Maguire, et al. Has the earth’s sixth mass extinction  
497 already arrived? *Nature*, 471(7336):51–57, 2011.
- 498 [4] Mickael Binois, Robert B Gramacy, and Mike Ludkovski. Practical het-  
499 eroscedastic gaussian process modeling for large simulation experiments. *Jour-  
500 nal of Computational and Graphical Statistics*, 27(4):808–821, 2018.
- 501 [5] M Boner, Th Sommer, C Erven, and Hilmar Förstel. Stable isotopes as a tool  
502 to trace back the origin of wood. In *Proceedings of the international workshop,  
503 Fingerprinting methods for the identification of timber origins, October*, pages  
504 8–9, 2007.
- 505 [6] et al. Borbas, E. Terra/modis temperature and water vapor profiles 5-min l2  
506 swath 5km, 2015. URL: [http://dx.doi.org/10.5067/MODIS/MOD07\\_L2](http://dx.doi.org/10.5067/MODIS/MOD07_L2).  
507 061.
- 508 [7] Gabriel J Bowen and Justin Revenaugh. Interpolating the isotopic composition  
509 of modern meteoric precipitation. *Water resources research*, 39(10), 2003.
- 510 [8] Gabriel J Bowen, Leonard I Wassenaar, and Keith A Hobson. Global applica-  
511 tion of stable hydrogen and oxygen isotopes to wildlife forensics. *Oecologia*,  
512 143(3):337–348, 2005.
- 513 [9] Federica Camin, Markus Boner, Luana Bontempo, Carsten Fauhl-Hassek, Si-  
514 mon D. Kelly, Janet Riedl, and Andreas Rossmann. Stable isotope techniques  
515 for verifying the declared geographical origin of food in legal cases. *Trends in  
516 Food Science & Technology*, 61:176–187, 2017.

- 517 [10] Federica Camin, Luana Bontempo, Matteo Perini, and Edi Piasentier. Stable  
518 isotope ratio analysis for assessing the authenticity of food of animal ori-  
519 gin. *Comprehensive Reviews in Food Science and Food Safety*, 15(5):868–877,  
520 2016.
- 521 [11] John F Casale, James R Ehleringer, David R Morello, and Michael J Lott.  
522 Isotopic fractionation of carbon and nitrogen during the illicit processing  
523 of cocaine and heroin in south america. *Journal of Forensic Science*,  
524 50(6):JFS2005077–7, 2005.
- 525 [12] Lesley A Chesson, Janet E Barnette, Gabriel J Bowen, Craig S Cook, Charles B  
526 Douthitt, John D Howa, Janet M Hurley, Helen W Kreuzer, Michael J Lott,  
527 Luiz A Martinelli, Shannon P O’Grady, David W Podlesak, Brett J Tripple,  
528 Luciano O Valenzuala, and Jason B West. Applying the principles of isotope  
529 analysis in plant and animal ecology to forensic science in the americas. *Oe-  
530 cologia*, 187(4):1077–1094, 2018.
- 531 [13] Andrea Contina, Sarah Magozzi, Hannah B Vander Zanden, Gabriel J Bowen,  
532 and Michael B Wunder. Optimizing stable isotope sampling design in terrestrial  
533 movement ecology research. *Methods in Ecology and Evolution*, 13(6):1237–  
534 1249, 2022.
- 535 [14] Alexandre Courtiol, François Rousset, Marie-Sophie Rohwäder, David X Soto,  
536 Linn S Lehnert, Christian C Voigt, Keith A Hobson, Leonard I Wassenaar,  
537 and Stephanie Kramer-Schadt. Isoscape computation and inference of spatial  
538 origins with mixed models using the r package isorix. In *Tracking animal  
539 migration with stable isotopes*, pages 207–236. Elsevier, 2019.
- 540 [15] Noel Cressie. *Statistics for spatial data*. John Wiley & Sons, 2015.
- 541 [16] Marine Cusa, Katie St John Glew, Clive Trueman, Stefano Mariani, Leah Buck-  
542 ley, Francis Neat, and Catherine Longo. A future for seafood point-of-origin  
543 testing using DNA and stable isotope signatures. *Reviews in Fish Biology and  
544 Fisheries*, 32(2):597–621, June 2022.
- 545 [17] Andreas Damianou and Neil D Lawrence. Deep gaussian processes. In *Artificial  
546 intelligence and statistics*, pages 207–215. PMLR, 2013.
- 547 [18] V. Deklerck. Timber origin verification using mass spectrometry: Challenges,  
548 opportunities, and way forward. *Forensic Science International: Animals and  
549 Environments*, 3:100057, 2023.
- 550 [19] Sandra Díaz, Josef Settele, Eduardo S Brondízio, Hien T Ngo, John Agard,  
551 Almut Arneth, Patricia Balvanera, Kate A Brauman, Stuart HM Butchart,  
552 Kai MA Chan, et al. Pervasive human-driven decline of life on earth points  
553 to the need for transformative change. *Science*, 366(6471):eaax3100, 2019.
- 554 [20] Eleanor E Dormontt, Markus Boner, Birgit Braun, Gerhard Breulmann, Bernd  
555 Degen, Edgard Espinoza, Shelley Gardner, Phil Guillery, John C Hermanson,

- 556 Gerald Koch, et al. Forensic timber identification: It's time to integrate disci-  
557 plines to combat illegal logging. *Biological Conservation*, 191:790–798, 2015.
- 558 [21] Jacob R Gardner, Geoff Pleiss, David Bindel, Kilian Q Weinberger, and An-  
559 drew Gordon Wilson. Gpytorch: Blackbox matrix-matrix gaussian process in-  
560 ference with gpu acceleration. In *Advances in Neural Information Processing*  
561 *Systems*, 2018.
- 562 [22] Peter E Gasson, Cady A Lancaster, Roger Young, Sara Redstone, Isabella A  
563 Miles-Bunch, Gareth Rees, R Philip Guillery, Meaghan Parker-Forney, and Eliz-  
564 abeth T Lebow. Worldforestid: Addressing the need for standardized wood  
565 reference collections to support authentication analysis technologies; a way  
566 forward for checking the origin and identity of traded timber. *Plants, People,*  
567 *Planet*, 3(2):130–141, 2021.
- 568 [23] Yuri Gori, Ana Stradiotti, and Federica Camin. Timber isoscapes. a case study  
569 in a mountain area in the italian alps. *PLoS One*, 13(2):e0192970, 2018.
- 570 [24] Carlos Guestrin, Andreas Krause, and Ajit Paul Singh. Near-optimal sensor  
571 placements in gaussian processes. In *Proceedings of the 22nd international*  
572 *conference on Machine learning*, pages 265–272, 2005.
- 573 [25] Robert J. Hijmans. terra: Spatial data analysis, 2022. R package version  
574 1.6-17. URL: <https://CRAN.R-project.org/package=terra>.
- 575 [26] Micha Horacek, Michael Jakusch, and Hannes Krehan. Control of origin of larch  
576 wood: discrimination between european (austrian) and siberian origin by stable  
577 isotope analysis. *Rapid Communications in Mass Spectrometry*, 23:3688–3692,  
578 2009.
- 579 [27] Huffman, G.J. and Behrangi, A. and Bolvin, D.T. and Nelkin, E.J. Gpcp  
580 version 3.1 satellite-gauge (sg) combined precipitation data set, 2020. URL:  
581 [https://disc.gsfc.nasa.gov/datasets/GPCPMON\\_3.1/summary](https://disc.gsfc.nasa.gov/datasets/GPCPMON_3.1/summary).
- 582 [28] Martin Jankowiak, Geoff Pleiss, and Jacob Gardner. Parametric gaussian pro-  
583 cess regressors. In *International Conference on Machine Learning*, pages 4702–  
584 4712. PMLR, 2020.
- 585 [29] Akira Kagawa and Steven W Leavitt. Stable carbon isotopes of tree rings as  
586 a tool to pinpoint the geographic origin of timber. *Journal of Wood Science*,  
587 56(3):175–183, 2010.
- 588 [30] Kaushi ST Kanankege, Moh A Alkhamis, Nicholas BD Phelps, and Andres M  
589 Perez. A probability co-kriging model to account for reporting bias and recog-  
590 nize areas at high risk for zebra mussels and eurasian watermilfoil invasions in  
591 minnesota. *Frontiers in veterinary science*, 4:231, 2018.
- 592 [31] Jason Kirk, Joaquin Ruiz, John Chesley, Spencer Titley, and Spence Titley.  
593 The origin of gold in south africa: Ancient rivers filled with gold, a spectacu-  
594 lar upwelling of magma and a colossal meteor impact combined to make the

- 595 witwatersrand basin a very special place. *American Scientist*, 91(6):534–541,  
596 2003.
- 597 [32] Geoff Koehler, Kevin J Kardynal, and Keith A Hobson. Geographical assign-  
598 ment of polar bears using multi-element isoscapes. *Scientific Reports*, 9(1):1–9,  
599 2019.
- 600 [33] Andreas Krause and Carlos Guestrin. Nonmyopic active learning of gaussian  
601 processes: an exploration-exploitation approach. In *Proceedings of the 24th*  
602 *international conference on Machine learning*, pages 449–456, 2007.
- 603 [34] Kailin Kroetz, Gloria M Luque, Jessica A Gephart, Sunny L Jardine, Patrick  
604 Lee, Katrina Chicojay Moore, Cassandra Cole, Andrew Steinkruger, and C Josh  
605 Donlan. Consequences of seafood mislabeling for marine populations and  
606 fisheries management. *Proceedings of the National Academy of Sciences*,  
607 117(48):30318–30323, 2020.
- 608 [35] Naoki Kurashima, Yukiko Makino, Setsuko Sekita, Yasuteru Urano, and Tet-  
609 suo Nagano. Determination of origin of ephedrine used as precursor for illicit  
610 methamphetamine by carbon and nitrogen stable isotope ratio analysis. *Ana-*  
611 *lytical chemistry*, 76(14):4233–4236, 2004.
- 612 [36] Jin Li and Andrew D Heap. A review of spatial interpolation methods for  
613 environmental scientists. 2008.
- 614 [37] Chao Ma, Hannah B Vander Zanden, Michael B Wunder, and Gabriel J Bowen.  
615 assignr: an r package for isotope-based geographic assignment. *Methods in*  
616 *Ecology and Evolution*, 11(8):996–1001, 2020.
- 617 [38] Naren Ramakrishnan, Chris Bailey-Kellogg, Satish Tadepalli, and Varun N  
618 Pandey. Gaussian processes for active data mining of spatial aggregates. In  
619 *Proceedings of the 2005 SIAM International Conference on Data Mining*, pages  
620 427–438. SIAM, 2005.
- 621 [39] Saeida Saadat, Hardi Pandya, Aayush Dey, and Deepak Rawtani. Food foren-  
622 sics: Techniques for authenticity determination of food products. *Forensic*  
623 *Science International*, page 111243, 2022.
- 624 [40] Nele Schmitz, Volker Haag, Céline Blanc-Jolivet, Markus Boner, María Teresa  
625 Cervera, Manuel Chavesta, Richard Cronn, Victor Deklerck, Carmen Diaz-  
626 Sala, Eleanor Dormontt, Peter Gasson, David Gehl, John C. Hermanson,  
627 Eurídice Honorio Coronado, Cady Lancaster, Frederic Lens, Estephanie Pa-  
628 tricia Liendo Hoyos, Sandra Martínez-Jarquín, Rolando Antonio Montenegro,  
629 Kathelyn Paredes Villanueva, Tereza Cristina Monteiro Pastore, Tahiana Ra-  
630 mananantoandro, Harisoa Ravaomanalina, Alexandre Magno Sebbenn, Niklas  
631 Tysklind, Mart Vlam, Charlie Watkinson, and Michael Wiemann. General sam-  
632 pling guide for timber tracking. How to collect reference samples for timber  
633 identification. *Global Timber Tracking Network*, 2019.

- 634 [41] Rolf TW Siegwolf, J Renée Brooks, John Roden, and Matthias Saurer. *Stable isotopes in tree rings: inferring physiological, climatic and environmental responses*. Springer Nature, 2022.
- 635  
636
- 637 [42] Anthony J. Silva, Rosalee S. Hellberg, and Robert H. Hanner. Chapter 7 -  
638 seafood fraud. In Rosalee S. Hellberg, Karen Everstine, and Steven A. Sklare,  
639 editors, *Food Fraud*, pages 109–137. Elsevier, 2021.
- 640 [43] Katie St. John Glew, Laura J Graham, Rona AR McGill, and Clive N Trueman.  
641 Spatial models of carbon, nitrogen and sulphur stable isotope distributions  
642 (isoscapes) across a shelf sea: An inla approach. *Methods in Ecology and*  
643 *Evolution*, 10(4):518–531, 2019.
- 644 [44] Chang Su and Sargur Srihari. Evaluation of rarity of fingerprints in forensics.  
645 *Advances in Neural Information Processing Systems*, 23, 2010.
- 646 [45] H Swofford and C Champod. Probabilistic reporting and algorithms in forensic  
647 science: stakeholder perspectives within the american criminal justice system.  
648 *Forensic Science International: Synergy*, 4:100220, 2022.
- 649 [46] Nikolaas J Van der Merwe, JA Lee-Thorp, JF Thackeray, A Hall-Martin,  
650 FJ Kruger, H Coetzee, RHV Bell, and M Lindeque. Source-area determination  
651 of elephant ivory by isotopic analysis. *Nature*, 346(6286):744–746, 1990.
- 652 [47] Peter van der Sleen, Pieter A Zuidema, and Thijs L Pons. Stable isotopes  
653 in tropical tree rings: theory, methods and applications. *Functional Ecology*,  
654 31(9):1674–1689, 2017.
- 655 [48] Charles J Watkinson, Peter Gasson, Gareth O Rees, and Markus Boner. The  
656 development and use of isoscapes to determine the geographical origin of quer-  
657 cus spp. in the united states. *Forests*, 11(8):862, 2020.
- 658 [49] Charles J Watkinson, Gareth O Rees, Sabine Hofem, Lina Michely, Peter Gas-  
659 son, and Markus Boner. A case study to establish a basis for evaluating geo-  
660 graphic origin claims of timber from the solomon islands using stable isotope  
661 ratio analysis. *Frontiers in Forests and Global Change*, 4, 2022.
- 662 [50] Jason B West, Gabriel J Bowen, Todd E Dawson, and Kevin P Tu. *Isoscapes:  
663 understanding movement, pattern, and process on Earth through isotope map-  
664 ping*. Springer, 2010.
- 665 [51] Christopher KI Williams and Carl Edward Rasmussen. *Gaussian processes for  
666 machine learning*, volume 2. MIT press Cambridge, MA, 2006.
- 667 [52] Barry T Wilson, Andrew J Lister, Rachel I Riemann, and Douglas M Griffith.  
668 Live tree species basal area of the contiguous united states (2000-2009). 2013.
- 669 [53] Michael B Wunder. Using isoscapes to model probability surfaces for deter-  
670 mining geographic origins. *Isoscapes: understanding movement, pattern, and  
671 process on Earth through isotope mapping*, pages 251–270, 2010.

672 [54] Stefan Ziegler, Stefan Merker, Bruno Streit, Markus Boner, and Dorrit E Jacob.  
673 Towards understanding isotope variability in elephant ivory to establish isotopic  
674 profiling and source-area determination. *Biological Conservation*, 197:154–163,  
675 2016.

## Supplementary material

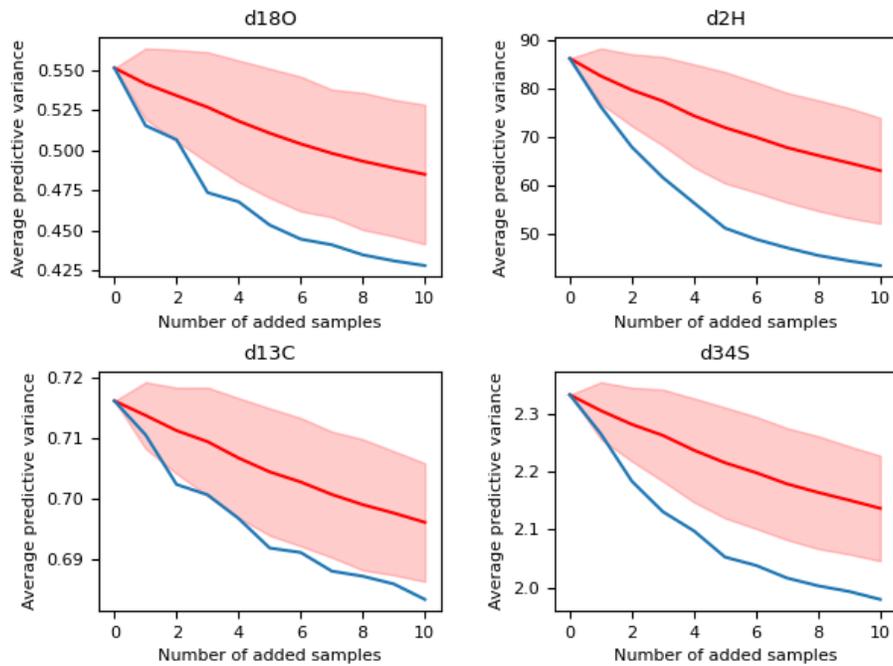


Figure 5: Average predictive variances for  $\delta^{18}O, \delta^2H, \delta^{13}C$  and  $\delta^{34}S$  as a function of the number of samples added to the base training data set; blue - active learning strategy; red - random sampling (shaded area denotes values within two standard deviations of the mean across  $n_r = 100$  simulations).