

1 **Authors:**

2 Corresponding author:

3 **Quinn Asena:**

- 4 • qasena@wisc.edu
- 5 • Address from which the work was done:
- 6 School of Environment, University of Auckland, 23 Symonds Street, Auckland, New Zealand
- 7 1010.
- 8 • Present address: Department of Geography, University of Wisconsin-Madison, Science Hall,
- 9 Madison, WI, 53703

10 Co-authors:

11 **George L. W. Perry:**

- 12 • george.perry@auckland.ac.nz
- 13 • School of Environment, University of Auckland, 23 Symonds Street, Auckland, New Zealand
- 14 1010.

15 **Janet M. Wilmshurst:**

- 16 • Wilmshurstj@landcareresearch.co.nz
- 17 • [Research](#)-Manaaki Whenua – Landcare Research, 76 Gerald Street, Lincoln, New Zealand,
- 18 7608

19

20 **Keywords:** palaeoecology; paleolimnology; ecological modelling; pseudoproxy; proxy
21 system modelling; virtual ecology; uncertainty; pseudoproxy modelling;

Information loss in palaeoecological data from process and observer error

Quinn Asena, George L. W. Perry, and Janet M. Wilmshurst

Abstract

Palaeoecological data ~~give us~~provide insight into how ecosystems have changed in the past, and, with the development of new sources of proxy data and statistical methods, they are being used to address questions around the underlying mechanisms of change, such as biotic- and climate-ecosystem interactions. However, inferences from palaeoecological data can be hindered by uncertainties inherent in core-type samples that arise from environmental processes and observer-introduced error. Environmental processes, core extraction methods, sub-sampling strategies, laboratory methods, and data processing can potentially mask 'true' signals in the data. ~~The influence of~~How different sources of uncertainty ~~on~~influence the inferences drawn from palaeoecological data ~~are~~is rarely assessed, but ~~are~~is critical to the confidence of our conclusions. To address this concern, we use a virtual ecological approach to assess the ~~influences~~effects of environmental and observer introduced uncertainty to better understand which of them ~~have the~~strongestmost influence ~~on~~statistical methods applied to the data. Quantifying information loss from uncertainty can ~~be used to~~inform study design before a project is carried out ~~to~~, and so increase the likelihood of detecting a given signal of interest and make more robust inferences from statistical analyses of palaeoproxy data. We generate synthetic 'error-free' core-type samples of pseudoproxies, on which environmental and observational processes are systematically introduced to impose uncertainties on the simulated pseudoproxies. The influence of three sources of uncertainty (core mixing, sub-sampling, and proxy quantification from sub-samples), are assessed for their individual and combined effects on two statistical methods ~~used to synthesise palaeoecological records~~: Fisher Information and principal curves. Increasing sub-sampling intervals has the most ~~substantial~~ influence on the two statistical methods applied to the pseudoproxy data. When combined, the interaction between increasing sub-sampling interval, and decreasing the number of proxies counted per sub-sample has the strongest influence on Fisher Information and principal curves. Fisher Information and principal curves are not affected in the same way by introducing uncertainty, with principal curves being less influenced by simulated proxy counting and sub-sampling of the core. Virtually assessing uncertainties is a powerful method to better understand the influence that uncertainties introduced at different parts of the analytical process have on conclusions drawn from palaeoecological data.

1 Introduction

Palaeoecological data extends the temporal extent over which we can investigate ecosystem change well beyond the observational record (Kosnik and Kowalewski, 2016). These long-term records are crucial for understanding ecosystem trajectories and climate-ecosystem interactions, as such dynamics may unfold over centuries or millennia (Jackson,

2007). ~~A growing number of~~ Many proxies and statistical methods are ~~available~~ used to address how ecosystems respond to environmental and human pressures through time. Palaeoecologists ~~seek to~~ use these data to go beyond describing past changes in, for example, the relative abundances of species, to uncover underlying mechanisms of change such as biotic interactions and species-environment relationships (Williams, ~~Blois, and Shuman et al.~~, 2011). However, the inferences drawn from palaeoecological data may be ~~hindered or even~~ limited by their uncertainties, such as a paucity of observations over time and space, environmental degradation of samples, and observer-introduced error. Thus, to make robust inferences from palaeoecological data, such uncertainties need to be better understood and quantified. [Here, we focus on proxies reflecting species community data such as fossil pollen extracted from a core sample and quantified on microscope slides.](#)

1.1 Uncertainties in palaeoecological data

Palaeoecological data derived from core-type samples are subject to numerous uncertainties, including: (i) environmental effects and landscape processes affecting the sample before extraction; (ii) the methods used in the field to extract the sample; (iii) laboratory techniques applied to extract and quantify data from the core; and (iv) quantitative analyses applied to the data (Table 1). Environmental processes and observation error can affect the representation of species in the data (e.g., their observed relative abundances; Goring et al., 2013). Sub-sampling strategies may alter the chronological placement of events (Liu et al., 2012; Parnell et al., 2008). Manipulation of data, such as interpolation to satisfy statistical assumptions, can introduce statistical artifacts and increase type-I error rates: [\(i.e., a false positive\)](#). Such uncertainties affect the robustness of statistical ~~methods~~ results and the inferences we may draw from them.

Table 1: Sources of uncertainty from pre-sampling natural processes to statistical analyses of data. Uncertainties are not independent across categories and can propagate through the observational process and subsequent analyses.

Source of Uncertainty	Examples
Physical, chemical and biological processes acting on the core or proxy. Not introduced by the observer.	Hiatuses, catchment erosion, variable sedimentation rates and mixing, bioturbation, changing sources of sediment or peat over time, preservation and taphonomy, occurrence of proxy in sample vs. actual abundance, differential preservation of proxies.
Observer introduced error from sampling collection and protocol.	Core compression during extraction, coring location (within basin or broader geographical context), sample depth/length and replication, core overlap, contamination.
Post-collection methods applied a sample and subsamples.	Contiguous/non-contiguous sub-sampling, sub-sampling resolution/density/thickness, sampling error/noise, proxy selection, taxonomic resolution, count method and proxy specific method error,

Data processing and interpretation.

dating frequency, dating precision and accuracy, observational error in proxy count.

Age-depth modelling, radiocarbon calibration, detrending, time-averaging, statistical methods selection, and understanding of proxy responses to environmental drivers.

85 1.2 Pseudoproxy experiments and virtual ecology

86 ~~One method that we~~ Virtual ecology (VE) can ~~use~~ be used to assess the influence of
87 uncertainty on statistical methods and associated inferences ~~is virtual ecology (VE)~~ (Zurell
88 et al., 2010). In the VE approach, simulated data are used as testbeds for recreating, in
89 simulation, observational processes. The synthetic data ~~act as~~ provide an 'error-free'
90 benchmark against which to assess simulated observational processes and analytical
91 methods. The underlying concept is that the synthetic data mimic the statistical properties
92 of empirical data without being subject to the same issues of, for example, limited grain size
93 or extent (Smerdon, 2012). Similarly, the simulated observational process aims to recreate
94 the statistical properties of the observer, such as the chance of observing a species
95 occurring at low abundances. Here, we adopt a ~~virtual ecological (VE)~~ VE approach to (i)
96 assess uncertainties introduced by environmental processes and observer-introduced
97 error by simulating data analogous to a sediment core-type sample, and (ii) to virtually
98 recreate environmental processes acting on a core, and the observational methods used to
99 extract data from the core sample. The VE approach follows the form: generate data →
100 simulate the observational process → analyse the 'true' and 'observed' data → assess the
101 analyses of the observed data against the 'true' data. We extend this approach to include
102 simulated environmental process errors that occur before the observational process. The
103 steps of our adapted VE approach: data generation, degradation, observation, analysis, and
104 assessment are described in sections 2.1 through 2.3. Empirical data typically lack an
105 'error-free' control, and even high-quality empirical data (e.g., highly resolved proxy data
106 from a laminated lake sediment core) incorporate multiple uncontrollable sources of
107 uncertainty. Virtual experimentation allows for the effects of sources of uncertainty to be
108 explored for their individual and interaction effects in a systematic and controlled way
109 (Smerdon, 2012).

110 The VE approach is similar to pseudoproxy experiments where modified observational
111 data or pseudoproxies (i.e., simulated proxy data ('pseudoproxies')) are used in place of
112 empirical observations (Mann and Rutherford, 2002) and analysed in the same way as
113 empirical measurements (e.g. Asena, Perry and Wilmshurst, et al., 2024). Pseudoproxy
114 experiments ~~originate from~~ originated in climatology where they are a method of assessing
115 palaeoclimate reconstructions (Mann and Rutherford, 2002; Christiansen, Schmith, and
116 Thejll et al., 2009; Bothe, Wagner, and Zorita et al., 2019) and,]; here, we apply the same
117 concepts to palaeoecological data. We use the term virtual ecology to describe the approach
118 by which synthetic data are generated, sampled from and analysed in ways comparable to
119 empirical data (Zurell et al. 2010). ~~The term 'pseudoproxy' is used to refer to the simulated~~
120 ~~proxy data themselves., 2010).~~ While simulated data cannot substitute entirely for reality,

121 they provide an experimental platform (*sensu* Peck, 2004) with which to ~~better~~ understand
122 the processes that influence the formation and analysis of empirical data.

123 Pseudoproxies have been widely used in climatology (Mann and Rutherford, 2002; Bothe,
124 ~~Wagner, and Zorita et al.~~, 2019), but virtually assessing sampling methods and statistical
125 approaches on palaeoecologically relevant data is less common (although see Asena,~~Perry~~
126 ~~and Wilmshurst, et al.~~, 2024; ~~Blaauw, Bennett, and Christen, Blaauw et al.~~, 2010; and Benito,
127 ~~Gil-Romera, and Birks, et al.~~, 2020). The lack of understanding ~~around of how~~ uncertainties
128 in palaeoecological data ~~has raised concerns regarding affect~~ the ~~reliability of~~ inferences
129 ~~drawn made~~ from them ~~has raised concern~~ (Blaauw 2012; Blaauw,~~Christen, and Aquino-~~
130 ~~López et al.~~, 2020). We address this knowledge gap by:

- 131 (i) generating multivariate pseudoproxy archives (considered analogous to a
132 ~~time series of proxy data from a~~ core sample);
- 133 (ii) introducing environmental uncertainty to the pseudoproxy archives via
134 simulated core mixing;
- 135 (iii) introducing process and observer error by virtually recreating the
136 observational processes of sub-sampling the core and quantifying proxies
137 from the sub-samples; and
- 138 (iv) applying two multivariate statistical methods independently, Fisher's
139 Information (FI; ~~Fisher, 1922~~) and principal curves (PrC; ~~Hastie and Stuetzle,~~
140 ~~1989~~), to analyse the 'error-free' and degraded/sub-sampled data.
- 141 (v) ~~Featurefeature~~ analysis methods (a dimensional-reduction method that
142 collapses time-series to a set of metrics) are then applied to ~~each~~ the FI and
143 PrCs ~~separately~~ to quantify the effect of increasing levels of uncertainty ~~on~~
144 ~~the two example methods.~~

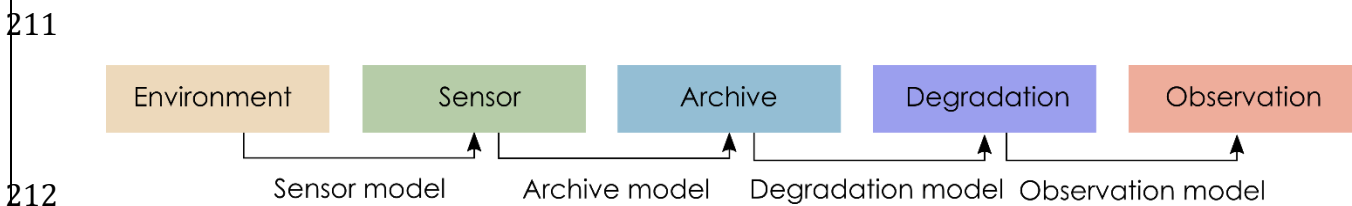
145 Our overarching aim is to quantify the information loss in palaeoecological analyses from
146 environmental uncertainties, and process and observer error, and how this influences
147 statistical analysis and inference using such data. We use ~~PrCsFI~~ and ~~FIPrCs~~ as examples
148 that capture the underlying drivers of a system in different ways. ~~PrC'sFI captures shorter-~~
149 ~~term variance e.g., those driven by short-term stochastic processes. PrCs,~~ as a method of
150 indirect gradient analysis, primarily captures the long-term ecosystem trends in the
151 pseudoproxies resulting from the primary driver in the scenarios. ~~FI captures shorter-term~~
152 ~~variance driven by stochastic processes such as the random walk driver.~~ A better
153 understanding of how individual and combined sources of uncertainty affect statistical
154 results and interpretation can help inform study design (e.g., determining the number of
155 replicate cores or the sub-sampling resolution needed to detect a signal of interest) and the
156 confidence in the statistical results of a study.

157

196 2 Methods

197 2.1 Simulating pseudoproxies

198 Pseudoproxies are simulated using the model described in Asena, Perry, and Wilmshurst et
199 al. (2024), following the proxy system model (PSM) framework (Evans et al. 2013) where a
200 sensor (e.g., terrestrial vegetation) responds to environmental drivers and records that
201 response via proxies (e.g., fossil pollen) counted in an archive such as a lake sediment (in
202 this case a pseudoproxy record). The model we use follows the conceptual framework of a
203 PSM, but is not process-based. We adapt the PSM conceptual framework (Evans et al.,
204 2013) to explicitly include degradation models that describe archive-altering processes
205 before observations are drawn (Figure 1). Asena et al., (2024) describe the three sub-
206 models that generate the pseudoproxy data: (i) a driver model representing the
207 environment in which the archive forms; (ii) a sensor model representing the response of a
208 sensor (e.g., terrestrial vegetation) to the environment; and (iii) an archive model that
209 represents how the response of the sensor is recorded (e.g., as fossil pollen) in a medium
210 such as a sediment core.



211
212
213 *Figure 1: Adapted Proxy System Model framework (Evans et al., 2013) describing the process*
214 *by which environmental drivers act on a sensor, which records its response in an archive, from*
215 *which observations are drawn. We have included degradation models to describe processes*
216 *acting on the archive before observations are drawn.*

217 'Error-free' pseudo proxies are simulated - in the environment, sensor, and archive sub-
218 models (Figure 1). The model consists of four components: (i) environmental driver change
219 over time; (*environment*); (ii) species' niches with respect to the driving environment (the
220 *sensor* response); (iii) pseudoproxy abundances recording the response to driver change in
221 an archive; (*archive record*); and (iv) a representation of the formation of the core. (*archive*
222 *characteristics such as core length and accumulation rate*). In summary, the model
223 generates archives of pseudoproxies consisting of 200 potential species representing a
224 palaeoecological record free from the process and observer error associated with empirical
225 data. Pseudoproxy abundances are simulated as a response to extrinsic and dynamic
226 environmental drivers with intrinsic variability from disturbance events, introduction to
227 the population via dispersal, and variation in carrying capacity over time. Each species has
228 a tolerance for each environmental driver that, together, defines the species niche and
229 determines the population growth rate of a species at any given time-step as a function of
230 the environmental drivers. If any of the environmental drivers fall outside of a species
231 tolerance to that driver, the species will have a negative growth rate and may eventually
232 become locally extinct. Species that are tolerant of the current environmental conditions
233 can be introduced via dispersal, thus creating a species turnover as conditions change.

234 Simulating 200 species covers a wide range of the driver parameter-space and allows
235 different species assemblages to emerge as driver conditions change. Only a subset of the
236 200 species ~~exists within their niche (i.e., favourable driver conditions)~~ will be presented in
237 the simulated community at any one point in time.

238 Thirty-one replicate models were run for a duration of 5000 time-steps with a burn-in
239 period of 500 time-steps applied to allow species to stabilise with respect to the driving
240 environment. ~~The scenario we analyse~~ We aimed for a minimum of 30 replicates to account
241 for model variance, and 5000 time-steps (c. 5000 years) was sufficiently long to represent
242 ecological turnover in the model ~~The scenario we present here~~ has two environmental
243 drivers: (i) an abrupt environmental driver switching between constant conditions; and (ii)
244 a random walk driver weighted to 0.15 of the total environmental effect. The magnitude of
245 change in the abrupt driver is insufficient to cause a complete species turnover, and
246 generalist species survive the shift in extrinsic conditions. The random walk driver may
247 amplify or dampen the effects of the abrupt driver, but in general is favourable to most
248 species in the system. The parameters for each species are randomised for each replicate
249 model run around the same baseline parameters (~~supplementary table 2~~. Table SI2). The
250 results of an additional three scenarios with different environmental drivers can be found
251 in the supplementary information. We present the abrupt change scenario for two reasons:
252 (i) there are empirical examples of such events (e.g., the termination of the Younger Dryas
253 and the Bølling-Allerød warming; Williams et al., 2011), and (ii) an abrupt shift in
254 community composition is more likely to be detected by statistical analyses. If the
255 statistical methods we assess perform poorly on the abrupt scenario, they are unlikely to
256 perform better on more gradual community turnover.

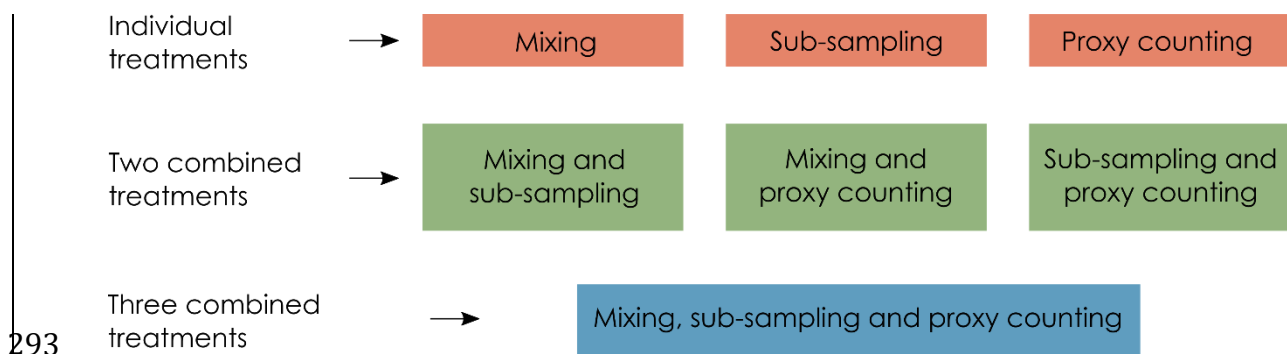
257

258 Each simulated core is characterised by simulated age, considered to be one year per time-
259 step, and an increase in depth per time-step. The accumulation rate is represented by a
260 combination of a linear decrease with time, with a smoothed random walk superimposed
261 to represent core compression in addition to landscape variability. Variable sedimentation
262 rates result in a different core length (and accumulation per time-step) for each replicate
263 simulation and change the number of model time-steps included in a sub-sample of one-
264 centimetre thickness. Variable change in depth per time-step is calculated as a smoothed
265 and scaled random walk (similar to Benito, ~~Gil-Romera, and Birks, et al.,~~ 2020)
266 representing landscape changes and possible hiatuses. A gradual decrease in depth with
267 age can occur from compaction or compression during extraction (Taranu et al., 2018).
268 The simulated data represent a core-type sample from which sub-samples are taken,
269 proxies quantified, and data analysed similar to real-world core samples. The simulated
270 core is used as an 'error-free' benchmark against which methodological and statistical
271 processes are assessed and ~~represent~~ represents the complete (and un-degraded) absolute
272 abundance of proxy data.

273 2.2 Degradation and sampling of pseudoproxies

274 ~~Asena, Perry, and Wilmshurst~~ In this section, we describe the (2024) ~~describe the three~~
275 ~~sub-models that generate the pseudoproxy data: (i) a driver model representing the~~

276 ~~environment in which the archive forms; (ii) a sensor model representing the response of a~~
 277 ~~sensor (e.g., terrestrial vegetation) to the environment; and (iii) an archive model that~~
 278 ~~represents how the response of the sensor is recorded (e.g., as fossil pollen) in a medium~~
 279 ~~such as a sediment core. In this section, we describe three more sub-models: (i) a~~
 280 degradation model representing post-depositional processes affecting the pseudoproxy
 281 data (in this case, core mixing); ~~(ii) an observational model), and the observation models~~
 282 [\(Figure 1\)](#) representing how proxies are quantified from a core sample (here, sub-sampling
 283 and the count method applied to [micro-fossils on microscope slides](#));. For details and ~~(iii)~~
 284 ~~a statistical model analysing both the ‘error-free’ data and degraded/sampled data. The~~
 285 ~~difference between the analysis visualisations of the ‘error-free’ data and~~
 286 ~~degraded/sampled data is then quantified. The pseudoproxy archive from degradation~~
 287 ~~models, see Asena et al., (2024). We apply three treatments to each of the 31-model~~
 288 replicates ~~has three treatments~~: mixing, [\(degradation\)](#), sub-sampling and proxy counting,
 289 [\(observation\)](#), applied at 10 levels each individually and combined [\(Figure 2\)](#). Together, the
 290 [degradation and observation models represent process- and observer- introduced error](#)
 291 [that affect the data before quantitative analyses are applied](#). Each of the 31 replicate
 292 archives results in 1210 datasets from the ‘error-free’ reference core to the most uncertain.



293
 294 [Figure 2: Each uncertainty \(treatment\) is applied to the ‘error-free’ pseudoproxy archive to](#)
 295 [systematically introduce uncertainty individually, and in combination.](#)

296

297 2.2.1 Virtual mixing: degradation model

298 Core mixing is applied as a centrally weighted rolling-average over time-windows of the
 299 archive ranging from unmixed (the ‘error-free’ benchmark) to a window of 10 time-steps.
 300 Simulated mixing is represented as consistent over time, rather than being a depth-
 301 dependent process such as bioturbation.

302 2.2.2 Virtual sub-sampling: observation model

303 An absolute proxy abundance per time-step is generated by the model and is sub-sampled
 304 at regular depth intervals ranging from one to ten centimetres. Each sub-sample has a
 305 thickness of one centimetre; the number of time-steps [withinspanned by](#) the sampled
 306 thickness is determined by the simulated accumulation rate. If the one-centimetre sub-
 307 sample covers multiple time-steps, the proxy abundances in that sub-sample are summed.

344 All sub-sample treatment data are converted to relative abundances before analysis so that
345 the analyses are not influenced by excessive values from the summed proxies of multiple
346 time-steps. When applied in combination with simulated proxy counting, values are
347 converted to relative abundance after the proxy count treatment.

348 2.2.3 Virtual proxy count: observation model

349 The process of counting proxies in a sub-sample, (e.g., counting [a few several](#) hundred
350 pollen grains or diatoms on a microscope slide,) is simulated by random sampling from the
351 absolute proxy abundances with resolutions increasing from 100 to 1000 by 100. The
352 probability of a proxy occurring in the random sample is based on the abundance of that
353 proxy. The sample is then converted to relative abundances comparable to empirical proxy
354 data. The simulated count treatment is applied to the raw, mixed, and sub-sampled data.
355 [Note that increasing](#)Increasing levels of uncertainty in the proxy counting treatment are
356 represented by a decrease in the proxy count resolution.

357 2.3 Quantitative analyses

358 [Fisher's information \(FI\)](#)After the application of the degradation and [principal curves](#)
359 [\(observation models, the pseudoproxy data span a gradient of uncertainty from the 'error-](#)
360 [free' data to the uncertain data that reflect what we observe from a system. We then apply](#)
361 [statistical analyses along this gradient to determine the influence of each source of](#)
362 [uncertainty \(the penultimate step of the VE approach\). FI and PrCs](#) are used to analyse
363 each treatment level from each replicate core. [Fisher's information \(treatment level being](#)
364 [the severity of each of the treatments, mixing, sub-sampling and proxy counting\).](#) FI and
365 ordination methods have been suggested as appropriate [analyses](#) where there are an
366 unlimited number of input variables of any data type that do not require *a priori*
367 knowledge of the driving state variables of a system (Roberts et al., 2018). [To quantify the](#)
368 [difference in the](#)The FI and PrC [among treatment levels produce a large amount of data, so](#)
369 [we use](#) feature analysis for time-series (FATs) [is used](#) (Nun et al., 2015).

370 [Because](#) to synthesise them and reduce each [time-series to a small number of metrics \(or](#) ←
371 ['features'\)](#) that we can compare across treatment levels (detailed below). Each replicate
372 core results in 1210 FI time-series and [principal curves, FATs is used to extract](#)PrCs, and by
373 [extracting a set of features from the FI time-series and the distance along the PrCs and](#)
374 [condense the information to a few features that are comparable among treatment levels.](#)
375 [The Euclidean distance between the extracted and scaled features is used, the difference](#)
376 [among treatment levels can be calculated](#) as a measure of distance [between treatment](#)
377 [levels, we use Euclidean distance.](#) The data analysis process is as follows: (i) calculate FI
378 and PrC for each treatment level; (ii) extract features from the FI and PrC outputs; and (iii)
379 calculate the distance between the features of each treatment level (Figure 1-3).

Format

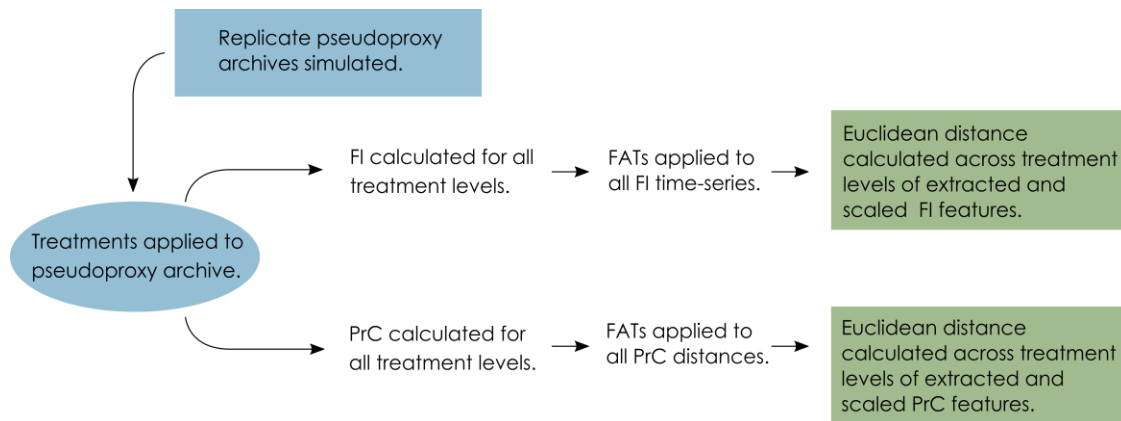


Figure 13: Conceptual data-flow of the degradation, sampling and analysis process. Treatments of mixing, sub-sampling and proxy counting are applied individually and in combination to the replicate pseudoproxy archives. (Figure 2). Fisher's information (FI) and principal curves (PrC) are applied separately to each treatment and subsequently analysed using feature analysis for time-series (FATs). Extracted features are scaled and Euclidean distance between each treatment level and the 'error-free' reference core is calculated.

2.3.1 Fisher Information (FI)

Fisher's information was developed as a method of quantifying to quantify information about an unknown parameter from measurable variables (Fisher, 1922) and has since been used as a measure of ecosystem stability in ecosystems (Cabezas and Fath, 2002; Mayer, Pawlowski, and Cabezas et al., 2006; Karunanithi et al. 2008; Eason et al., 2016). Fisher's information has been applied to palaeoecological data, with the suggestion that it can be used as an indicator of an approaching regime shift, such as a shift in diatom species' composition (Spanbauer et al., 2014). Fundamentally, FI evaluates the probability of detecting different system states over time, or the 'stability' of the system; thus, FI changes with the variance in the system. Pseudoproxies from each replicate core, and each level of uncertainty, are analysed using FI, using a custom R package (Asena, Young, and Pletzer et al., 2023).

2.4.3.2 Principal curves

Principal curves can be used to identify the underlying variables (e.g., an ecological gradient a steadily changing variable over time) that describe a system characterised by multiple state variables and can be considered as a form of non-linear principal components analysis (De'ath, 1999). In short, a PrC is a one-dimensional curve fit through the 'middle' of an n -dimensional space (e.g., species composition data; Hastie and Stuetzle, 1989) that can represent. The curve represents species composition compositions by mapping (or projecting) the sites, e.g., the species' composition at a given sample, onto a low-dimensional space, and using similarity or dissimilarity measures to measure the distance between sites (assessing, for example, the species' compositional change through time). The arrangement of sites reflects the composition of species in the reduced dimension space as the distance between sites is proportional to the distance in species

Formatt

Formatt

471 composition (De'ath, 1999). Principal curves can be used as a method of gradient analysis,
472 the underlying concept being that the species abundances change in a predictable way
473 along an ecological gradient. Here, we use PrCs to represent change in species composition
474 over time and as a method of indirect gradient analysis using the distances along the PrC.
475 Cubic smoothing splines are used to fit the PrC to the data. [Details for the implementation
476 of PrCs can be found in](#) Hastie and Stuetzle (1989), De'ath (1999), and Simpson and Birks
477 (2012). [detail the implementation of PrCs.](#) PrC analyses were conducted for all replicate
478 pseudoproxy datasets for each level of increasing uncertainty using the analogue package
479 (Simpson and Oksanen, 2020) in R (R Core Team 2020).

480 [2.4.1 Assessing results: Feature analysis for time-series](#)

481 [The final step of the VE approach is to assess the outcomes of the statistical analyses
482 applied to the degraded and sampled data against the 'error-free' pseudoproxies. To
483 compare the FI and PrC time-series across treatment levels, we use feature analysis.
484 Feature analysis is a method of reducing a two-dimensional time-series to a one-
485 dimensional set of 'features'. Features are metrics that describe a time-series in terms of
486 summary statistics \(e.g., mean and variance\), and more complicated descriptors such as
487 autocorrelation length \(Nun et al., 2015\). We developed a set of 62 features \(Table SI4\)
488 drawing on feature analysis for time-series \(Richards et al., 2011; Kim et al., 2011; Nun et
489 al., 2015; Sokolovsky et al., 2017\) and change point analysis \(Killick and Eckley, 2014\). The
490 individual and combined degradation and sampling treatments result in 1210 time-series
491 per replicate, yielding \$31 \times 1210 = 37510\$ virtual cores per scenario. ~~Thus, to compare the
492 analyses \(FI and PrCs\) of each pseudoproxy record, as uncertainties are introduced, with
493 the 'error-free' benchmark, we use feature analysis to reduce the dimensions of the time-
494 series for comparison across treatment levels.~~](#)

495 ~~Features are extracted from the FI time-series and the distances along the PrCs drawing on
496 feature analysis for time-series (FATs) (Richards et al. 2011; Kim et al. 2011; Nun et al.
497 2015; Sokolovsky et al. 2017) and change point analysis (Killick and Eckley 2014) to
498 describe the FI and PrC analyses as a series of metrics. Sixty-two features are extracted
499 from the FI time-series and PrC (supplementary table 4). Describing the FI time-series and
500 PrC as a series of metrics/features allows comparison between [the 1210](#) treatments by
501 calculating a distance measure (we use Euclidean distance) between the extracted features
502 of each treatment level. The feature analysis process is as follows:~~

- 503 1. Features are extracted from the FI time-series and distances along the PrCs for all
504 replicate model runs of the 'error-free' archive.
- 505 2. A correlation matrix of the features from the replicate 'error-free' datasets is
506 constructed, and features with ~~an absolute Pearson's~~ [Pearson's](#) correlation
507 coefficient greater than $|0.7|$ are excluded sequentially, recalculating the correlation
508 matrix until all highly correlated features are dropped, starting with the highest
509 correlation coefficient (Dormann, et al. 2007).
- 510 3. The remaining features (with ~~an absolute~~ [a](#) Pearson's correlation coefficient less
511 than $|0.7|$) are then calculated for all replicates across all treatment levels, resulting
512 in a number (ranging between 14-26 per scenario) of single metric features per
513 treatment that describe the FI time-series and PrC.

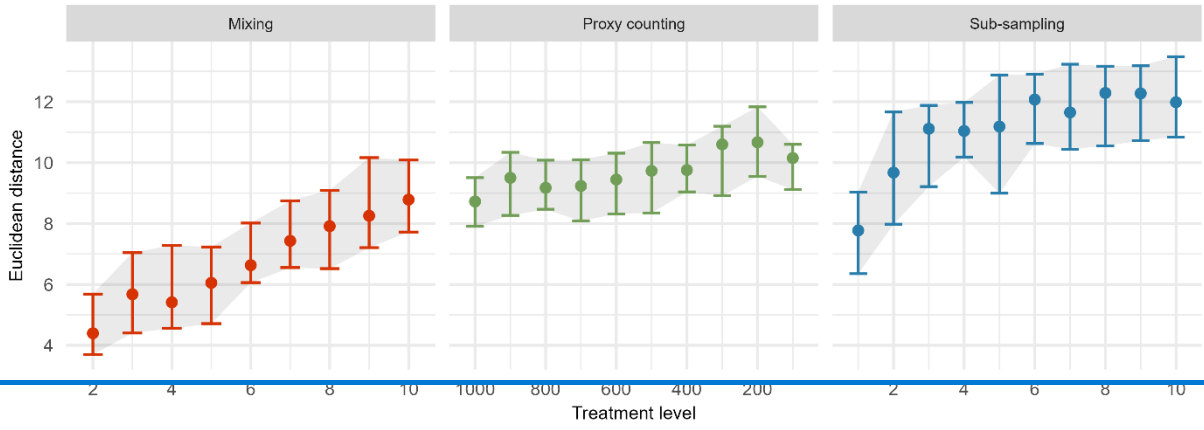
- 514 4. The features are scaled (by subtracting the mean of the entire series from each point
515 and dividing it by the series' standard deviation), and the Euclidean distance is
516 calculated across treatment levels, resulting in a single distance measure between
517 the 'error-free' core and each treatment level.
- 518 5. Summaries of the Euclidean distances are calculated for all treatment levels across
519 replicates, resulting in a single distance measure for each treatment from the 'error-
520 free' benchmark.
521

522 **3 Results**

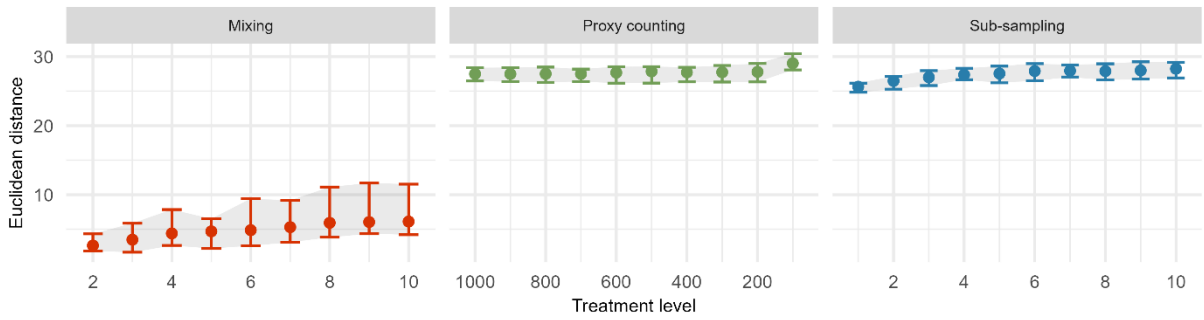
523 **3.1 Effects of individual sources of uncertainty**

524 In the extracted features for FI, sub-sampling causes the largest overall increase in the
525 median Euclidean distance from the 'error-free' core, followed by proxy counting and then
526 mixing; however, there is ~~some~~ overlap in the confidence envelopes across all three
527 treatments ~~((Figure Figure 24-A)~~. Distance ~~from the 'error-free' core~~ increases consistently
528 with uncertainty in the mixing treatment, showing a steeper increase in distance across
529 uncertainty compared with the other two treatments. Between successive treatment levels,
530 proxy counting shows little increase in the median Euclidean distance as uncertainty
531 increases. In the sub-sampling treatment, distance increases more in the lower uncertainty
532 levels than ~~in the~~ higher levels, ~~potentially~~ plateauing ~~somewhat~~ at higher uncertainty
533 ~~((Figure Figure 24-A)~~. ~~A degree of~~There is ~~some~~ variability ~~is visible~~ across the treatment
534 levels resulting from the stochasticity in the underlying model and simulated observational
535 processes (i.e., proxy counting is a random sampling process, and sub-sampling is
536 dependent on the variable accumulation rates of the core).

A - Fisher's information features



B - PrC distances features



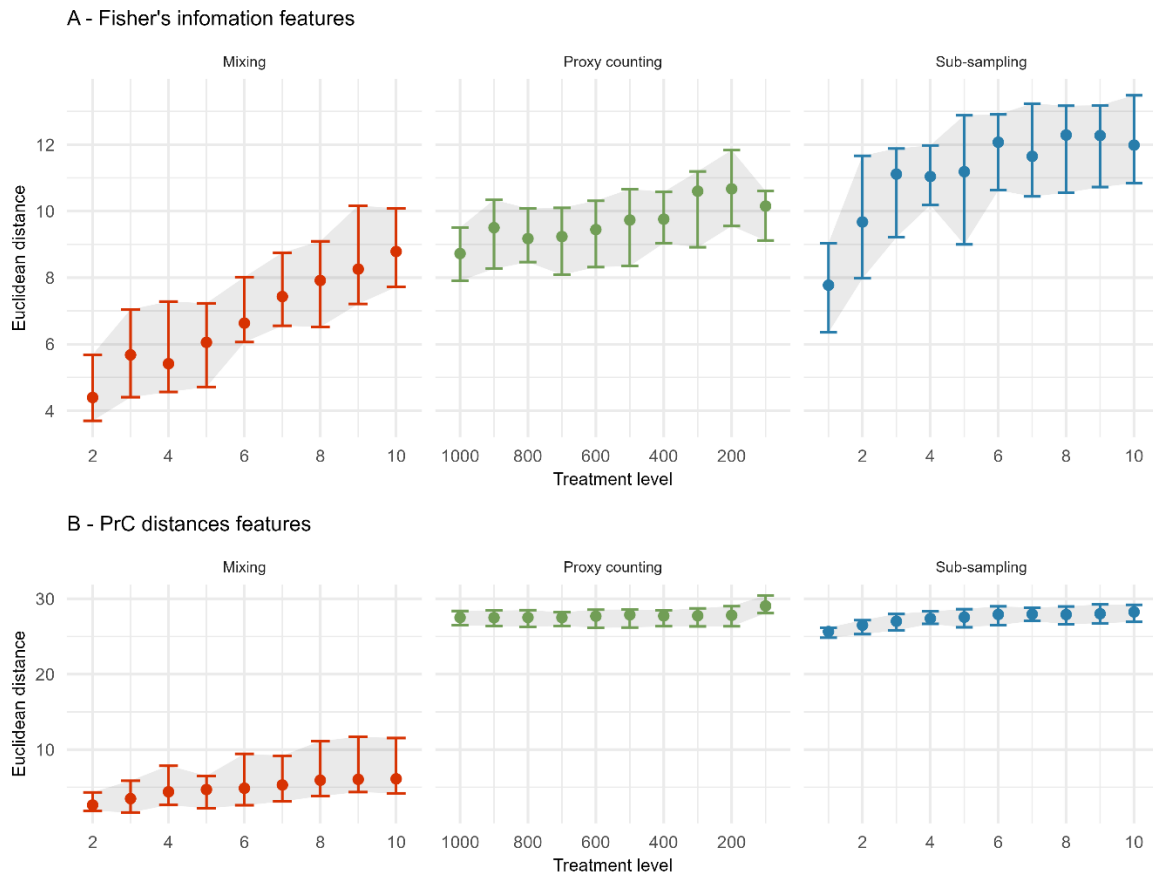


Figure 24: The median (dots), 25th, and 75th quantiles (error-bars and shaded area) of the Euclidean distance from the 'error-free' core of features extracted from the Fisher's information (A) and PrC distances (B) calculated across replicate simulations. Note the x-axis is organised so uncertainty consistently increases from left-to-right.

546 In the analysis of For the PrC features, across all scenarios the [mixing treatment has the](#)
 547 [least effect on median Euclidean distance is observed in the mixing treatment](#) ((Figure
 548 [Figure 24-B](#)). Proxy counting and sub-sampling have overlapping confidence envelopes,
 549 although they are much smaller than those of the mixing treatment. Proxy counting shows
 550 no consistent pattern in median Euclidean distance between successive treatment levels
 551 until the lowest count resolution. Conversely, the sub-sampling treatment shows an
 552 increase in the lower treatment levels, plateauing as sub-sampling interval increases
 553 ((Figure [Figure 24-B](#)).

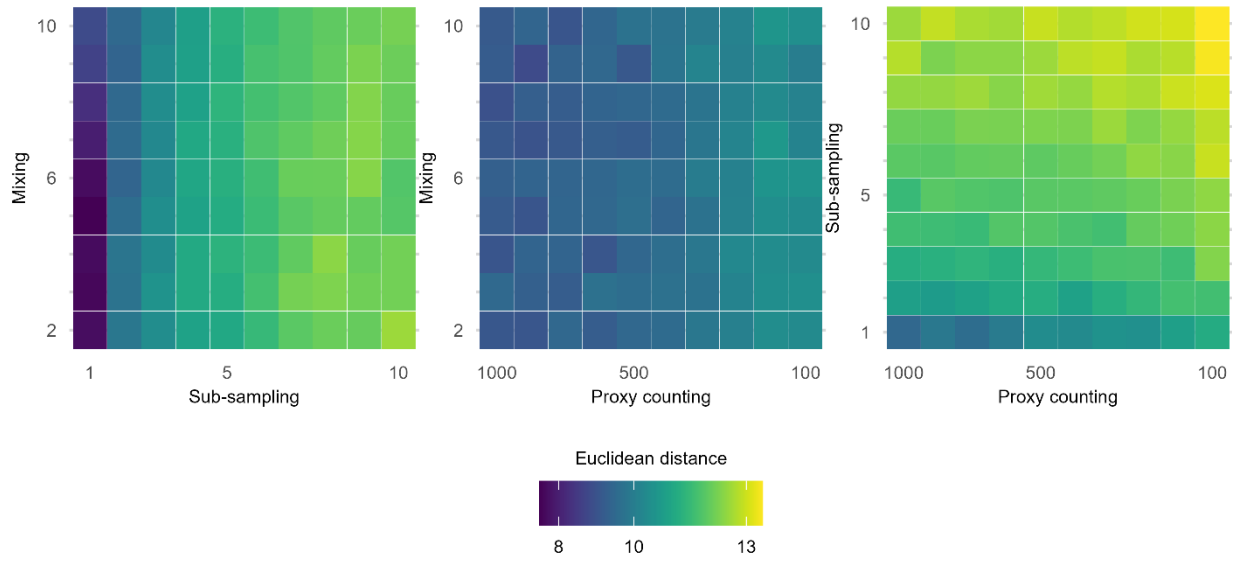
554

555 **3.2 Effects of two combined sources of uncertainty**

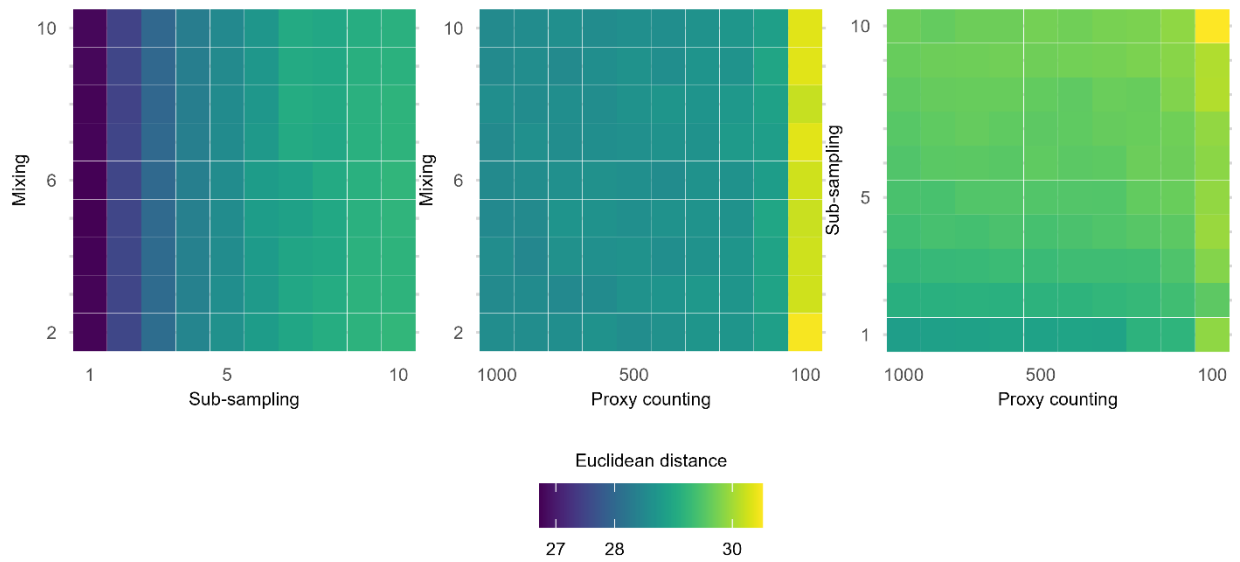
556 ~~Looking at individual treatment effects, the distance of the extracted features increases~~
557 ~~with the severity of the treatment levels.~~ In the following section, treatments are applied
558 simultaneously to determine which combinations cause the greatest effect on analyses of
559 the core. The greatest increase in mean Euclidean distance on the features extracted from
560 FI from the 'error-free' core arises from the interaction of the sub-sampling and proxy
561 counting uncertainties increasing together (i.e., along the diagonal; ~~Figure Figure 35-A~~),
562 with a stronger effect from the subsampling treatment. Mixing combined with sub-
563 sampling or with proxy counting, shows no clear interaction effect as the treatments
564 increase in severity together. The smallest increase in distance is from the combination of
565 mixing with proxy counting, suggesting that sub-sampling tends to have the greatest
566 influence of the three treatments (~~Figure Figure 35-A~~). The effect of the mixing treatment
567 on the mean Euclidean distance is small when compared to those of either sub-sampling or
568 proxy counting (~~Figure Figure 35-A~~). Variability across the surface of each plot
569 ~~emerges~~emerges from underlying model stochasticity, random sampling in the simulated
570 count method, and the variable accumulation rates of the replicate cores.

571

A - Fisher's information features



B - Principal curve features



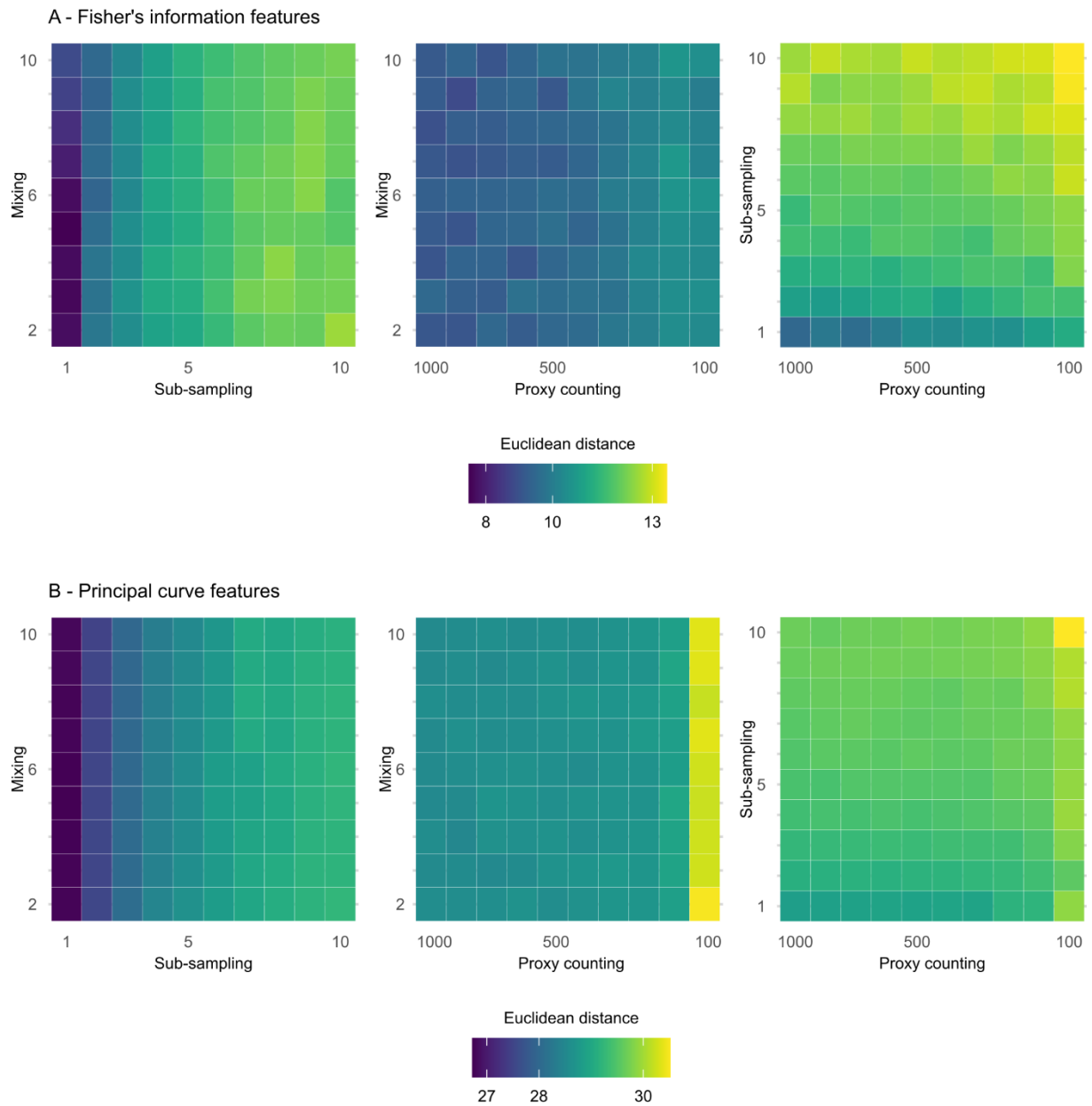


Figure 35: Mean Euclidean distance of features from the 'error-free' core of two treatments combined calculated across replicate simulations for Fisher's information (A) and principal curves (B). The mixing axis shows the number of time-steps over which mixing occurs. Along the sub-sampling axis, the frequency of sub-sampling in centimetres is shown, and the proxy counting axis displays count resolutions in number of individuals counted per sample. In the proxy counting treatment, uncertainty increases as count resolution decreases.

574 In the extracted features from the distances along the PrCs, the combined effects of sub-
 575 sampling with proxy counting show the largest increase in the mean Euclidean distance of
 576 all the combined treatments, with a weak interaction effect as proxy count and sub-
 577 sampling uncertainties increase together (Figure 35-B). No interaction effect is visible in
 578 either the combined treatments of mixing with sub-sampling or mixing with proxy
 579 counting (Figure 35-B).

580 3.3 Effects of three combined sources of uncertainty

581 ~~The interaction~~Interaction effects of all three uncertainties applied simultaneously ~~are were~~
 582 assessed for the extracted FI and PrC features ~~(Figure Figure 46)~~. An interaction effect is
 583 visible in the increase of mean Euclidean distance as the sub-sampling interval increases
 584 (along the x-axis), together with proxy counting resolution decreasing (across facets from
 585 top left to bottom right); however, no clear increase is visible along the mixing axis,
 586 indicating little contribution from mixing to the interaction of the three treatments
 587 (Figure 46). Overall, the greatest increase in the mean Euclidean distance among
 588 treatments is at the lowest proxy count and largest sub-sampling interval.

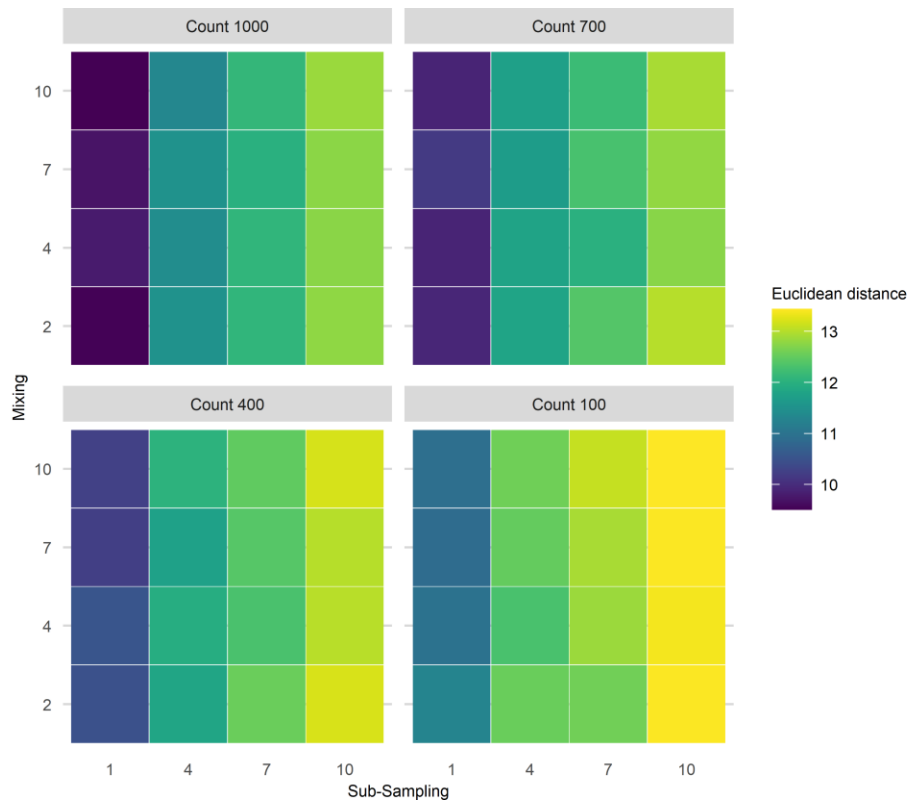


Figure 46: Mean Euclidean distance of Fisher's Information from the 'error-free' core for three treatments applied in combination. To demonstrate the three treatment dimensions, results are displayed such that each plot axis shows mixing (number of time-steps over which mixing occurs) and sub-sampling (frequency in centimetres) treatments, and each facet (sub-plot) is the proxy counting treatment (number of individuals counted in sub-sample). Uncertainty from proxy counting increases from the top left to the bottom right.

589 In the PrC features, applying three treatments in combination does not show a clear three-
 590 way interaction as uncertainty increases (Figure 3.47). An increase in mean Euclidean
 591 distance is visible along the sub-sampling axis (as sub-sampling interval increases), but
 592 there is little visible effect of the proxy counting treatment (reducing in resolution across
 593 facets from top left to bottom right) until the lowest resolution. Mixing contributes
 594 relatively little to the overall increase in mean Euclidean distance showing no visible
 595 increase along the mixing axis (Figure 3.47).

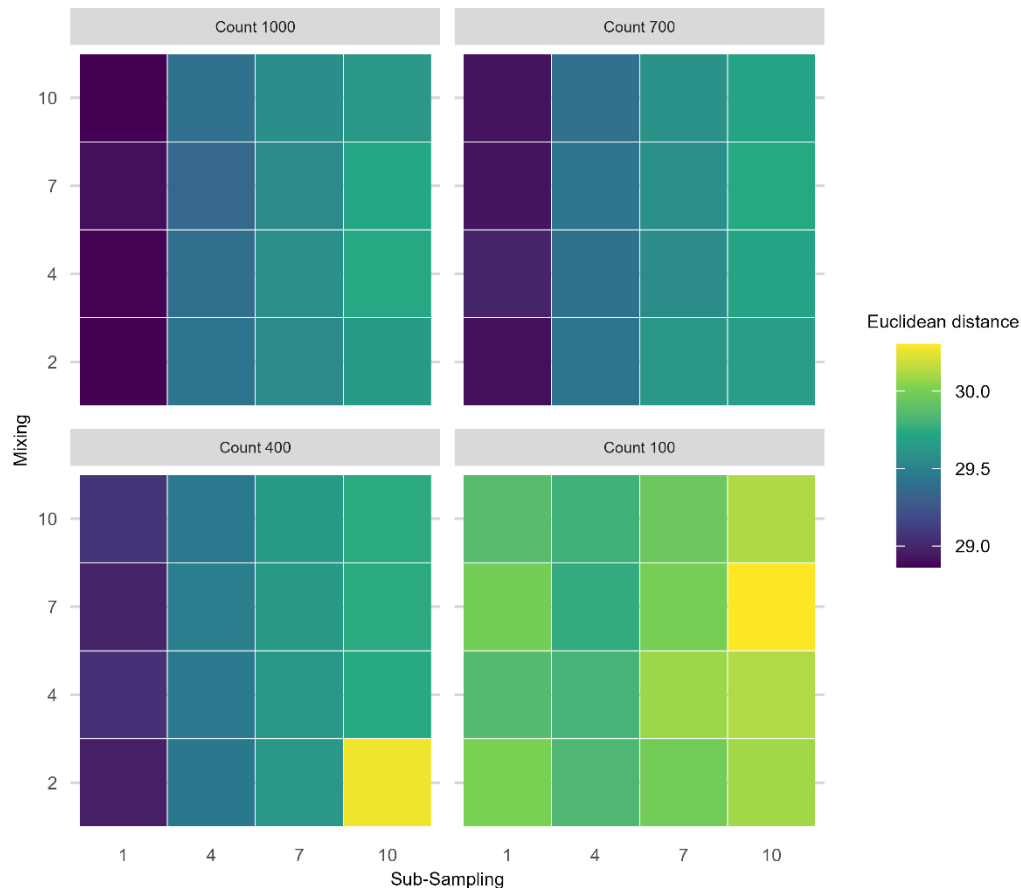


Figure 57: Mean Euclidean distance of the principal curves from the 'error-free' core of all uncertainties increasing in combination. Proxy count resolution decreases across facets from top left to bottom right. Within each facet the axes show the sub-sampling (frequency in centimetres) Same plot layout and mixing treatments (number of time-steps)-interpretation as Figure 6.

596 4 Discussion

597 To draw reliable conclusions from palaeoecological data, it is crucial to view our inferences
 598 in the context of the uncertainties they integrate. Our goal in this paper is Our goal in this
 599 paper was to address how process and observer error affect statistical methods applied to
 600 palaeoecological data and the inferences we draw from them. Here, we have assessed some
 601 of the uncertainties that affect species proxy data. Additionally, uncertainty arises from
 602 building chronologies (Blaauw et al., 2018; Parnell et al., 2008; Telford et al., 2004) and

603 measuring abiotic system variables, such as isotopic records and lake level reconstructions,
604 which carry their own uncertainties. ~~To draw reliable conclusions from palaeoecological~~
605 ~~data, it is crucial to view our inferences in the context of the uncertainties they integrate.~~

606 4.1 Effects of individual sources of uncertainty

607 Without exception, sub-sampling treatments show the largest effect on the median
608 Euclidean distances of the FI features (~~although,~~ with some overlap in the confidence
609 envelope with the proxy counting treatment), followed by the proxy counting process. For
610 the PrC features, both sub-sampling and proxy counting treatments ~~showhad~~ a similar
611 ~~magnitude of~~ effect on Euclidean distance. For both the FI and PrC features, the smallest
612 effect on Euclidean distance between the 'error-free' and degraded cores ~~iscame~~ from the
613 simulated mixing degradation process.

614 Perhaps the most interesting implication of our analyses in the context of empirical ecology
615 is how the effects of sources of error differ between the two statistical metrics (FI and PrC).
616 Beyond the initial increase in Euclidean distance in the PrC from the application of proxy
617 counting and sub-sampling (an unavoidable cost), there was little further effect until the
618 lowest proxy count and sub-sampling resolutions. Thus, PrCs may be a useful measure of a
619 system's trajectory for patchy data (e.g., initial analysis of low sub-sampling resolution data
620 before deciding where to focus sampling effort). For FI, treatment effects have a more
621 consistent increase in distance from the 'error-free' core with treatment intensity. The
622 short-term variability captured by FI may provide useful system indicators (*sensu* Eason
623 and Cabezas, 2012) but may also require high-quality data (e.g., high sub-sampling
624 resolution and proxy counting) for reliable inferences. The required temporal resolution of
625 the data is also likely to increase if accumulation rates are slow and species turnover is
626 rapid. Increased sub-sampling frequency is required in systems that change rapidly
627 compared with more stable systems where the difference between successive time-steps is
628 small. Slow accumulation rates mean that a one-centimetre-thick sub-sample integrates
629 multiple years of ecological change; thus, uncertainty from sub-sampling resolution will
630 increase if the accumulation rate is slow and ecological change is rapid. An observer may
631 interpret FI results with the knowledge of which sources of uncertainty, whether
632 controllable (e.g., sub-sampling and proxy counting resolution) or uncontrollable (e.g.,
633 mixing and accumulation rates), have the greatest influence. After introducing
634 uncertainties ~~to the pseudoproxies,~~ the primary patterns, such as long periods of increase
635 in FI, remain visible and offer a useful depiction of community change.

636 4.2 Effects of combined sources of uncertainty

637 For both FI and PrC, when two treatments are applied in combination, the greatest overall
638 increase in mean Euclidean distance from the 'error-free' core resulted from sub-sampling
639 uncertainty in combination with proxy counting uncertainty. ~~Although for~~ For the PrC the
640 maximum increase in mean Euclidean distance for two simultaneous treatments occurred
641 at the lowest proxy count resolution in combination with mixing, ~~but~~ such a low proxy
642 count resolution is unlikely in any empirical study. In our analyses, mixing has relatively
643 little effect compared with sub-sampling or proxy counting, and the effect tends to be
644 obscured by other sources of stochasticity (disturbance, dispersal, and temporal changes in

645 carrying capacity), variable accumulation rates, and randomised sampling of the proxy
646 abundances. [Similarly, when three sources of uncertainty are combined, mixing shows the](#)
647 [least effect, and no apparent interaction with sub-sampling or proxy counting. However, it](#)
648 [is worth noting, that our implementation of mixing is relatively simple in that the effect on](#)
649 [the pseudoproxy data is a centrally weighted smoothing. More extreme mixing, and time-](#)
650 [varied mixing effects are likely to show a stronger effect.](#)

651 In the case of FI, the greatest benefit in terms of reducing the distance from the ‘error-free’
652 core, [is derived](#) [came](#) from increasing sub-sampling resolution followed by proxy count
653 resolution. However, from the standpoint of an observer, increasing the resolution of an FI
654 time-series may not increase the information content, at least in terms of capturing long-
655 term patterns. Simple driving environments (e.g., the single driver simulations) are likely to
656 be adequately represented by FI applied to low-resolution data, [and; however,](#)
657 interpretation of FI from [complex driving](#) environments ([i.e., those](#) influenced by multiple
658 drivers) is potentially challenging without considerable knowledge of the underlying
659 driving conditions. The benefit of the VE approach is that it allows us to examine how the
660 underlying dynamics manifest in multivariate indicators of change, such as FI; however,
661 such near-complete information is rarely, if ever, available. Ultimately, an observer must
662 base [sampling and analysis](#) decisions, [such as sub-sampling strategy,](#) on their specific aims.
663 If accurate evaluation of short-term variability is a [requirementgoal](#) (e.g., for studies of
664 ecosystem resilience), then dedicating resources to sub-sampling resolution is likely to be
665 beneficial to analyses such as FI. Reducing observer error, a more controllable source of
666 error, may help detect a short-term signal of interest if the uncontrollable sources of error,
667 such as sediment accumulation rates, mixing, and driver variability, are sufficiently small
668 that the signal remains detectable.

669 PrC shows little increase in distance (after the initial increase from the application of the
670 treatments) from the ‘error-free’ reference with combined treatment levels. Thus, as a
671 representation of compositional change, it may be robust to low-resolution data (e.g.,
672 infrequent sub-sampling). Short-term changes in abundance (e.g., small disturbances) will
673 likely become less evident in the PrC; however, our results suggest that overall patterns
674 seen in a PrC are robust to multiple sources of error. Thus, from the perspective of an
675 observer, high-resolution data may only be required for PrC to identify short-term
676 compositional changes, such as perturbations that take a few generations to recover from.

677 **[4.3 Implications-Lessons from virtual ecology](#) for empirical studies**

678 [One purpose of using](#) [Using](#) virtual ecology to estimate the influence of sources of
679 uncertainty on quantitative analyses [ishelps us](#) to understand what can be done to mitigate
680 their effects and where to focus limited resources, such as time spent analysing an
681 individual core. Trade-offs are inherent in any sampling design. For example, is it more
682 advantageous to focus effort on spatial coverage by taking multiple cores rather than
683 increasing sub-sampling resolution on fewer cores? Such questions, of course, depend on
684 the intention of the study and knowledge of the study site/system ([i.e., some knowledge of](#)
685 [the uncertainties that will be encountered,](#) such as landscape changes through time).
686 Virtual ecology allows us to make a more informed decision about what field and
687 laboratory methods, and quantitative analyses will be most appropriate given the question

688 of interest. For example, if analysing a network of core data over a large geographic region,
689 where the observer is interested in the spatial consistency of the system's trajectory but
690 ~~doesthere are~~ not ~~have~~ the resources to extract highly resolved data from each core,
691 methods such as PrC may be more informative than FI. Conversely, if an observer is
692 interested in short-term change in a single core (or a region of particular interest in a
693 core), it may be worth allocating the time to extracting highly resolved data to increase the
694 reliability of analyses sensitive to variance such as FI. ~~Of course, FI and PrC are only two~~
695 ~~examples of a suite of available statistical analyses (Birks et al. 2012; Blaauw, Christen, and~~
696 ~~Aquino-López 2020) and an observer should apply more than one. However,~~ PrC and FI
697 provide different representations of a system's trajectory. Principal curves reflect the
698 system's overall trajectory (De'ath, 1999) and, as a form of indirect gradient analysis, PrC
699 reflects the primary driver of the scenario. In contrast, FI is more sensitive to short-term
700 variability (Eason et al., 2016) and so reflects driver interactions. Although FI does capture
701 the long-term system trajectory, these trends can be obscured by short-term variability
702 such as that caused by the random walk driver. Similar analyses (e.g., other ordination
703 methods) may be affected by sources of uncertainty in similar ways. ~~Of course, FI and PrC~~
704 ~~are only two of a suite of available statistical analyses (Birks et al. 2012; Blaauw et al.,~~
705 ~~2020), and an observer should apply more than one.~~

706 Alongside considerations of the sensitivity of analyses to uncertainty, ~~there~~ are questions of
707 how sensitive different proxies are to driver change and, consequently, how informative
708 the analyses are. The question of how different proxies respond to drivers at different
709 temporal and spatial scales (e.g., Wilmshurst et al., 2002) remains ~~largely~~
710 ~~unanswered, poorly resolved.~~ Interestingly, the sensitivity of different proxies to
711 environmental change may be ecosystem-specific. Phytoliths have been reported as more
712 sensitive than pollen to changes in dry forests, with the reverse being true of evergreen
713 forests at a site in Bolivia (Plumpton, Whitney, and Mayle et al., 2019). In savannahs, pollen
714 and phytoliths ~~have been shown to be~~ equally sensitive to changes in the environment
715 (Plumpton, Whitney, and Mayle et al., 2019). Thus, ecosystem-specific and proxy-specific
716 knowledge are important considerations, as increasing sub-sampling efforts to obtain a
717 higher resolution representation of change from numerical analyses may not be useful if
718 the proxy is not a reliable sensor at that resolution. Furthermore, compositional change is
719 not the sole (or necessarily the most appropriate) measure of system change, and other
720 ~~measures~~ ~~descriptors~~ such as body-size distributions, physiognomic, and functional and
721 phylogenetic diversity can be included (Goring et al., 2013; Reitalu et al., 2015; Clements
722 and Ozgul, 2016; Spanbauer et al., 2016; Adeleye et al., 2023).

723
724 The advantage of using a phenomenological modelling approach designed to mimic the
725 statistical properties of empirical data, is that the results are informative without
726 attempting to recreate species abundances of a specific system and addressing the
727 challenges faced by more process-based models. Our model and virtual approach can be
728 used to assess, for example, the performance of statistical methods, the influence of
729 observational and chronological uncertainty, and explore a range of driving environments.
730 It becomes possible to address questions such as 'what statistical method is likely to be
731 appropriate for my data and research question?' However, without recreating a specific

732 [system, it cannot answer questions such as ‘should I use a 2 or 4 cm resolution sub-](#)
733 [sampling procedure for this core?’.](#) To answer such questions, process-based VE
734 [approaches are necessary.](#) Of course, process-based models face numerous challenges of
735 [accurately representing mechanisms and being able to reasonably recreate a given system](#)
736 [to address such questions.](#) The VE approach can be applied using process-based models to
737 [generate data, or simulating data from a statistical model fit to the empirical data, and](#)
738 [following the same process of degrading, sub-sampling and fitting/re-fitting, to assess how](#)
739 [parameters change.](#)

740 [Finally, we can consider, empirical, experimental, semi-empirical approaches to](#) 741 [advance our understanding of uncertainties in palaeoecology.](#) ~~5-Conclusion~~

742 ~~Palaeoecological uncertainty can be considered at four levels: environmental processes,~~
743 ~~field methods, laboratory methods, and quantitative analyses (Table 1).~~ The effects of
744 ~~different sources of uncertainty are challenging to disentangle and quantify from empirical~~
745 ~~studies alone, and virtual ecology provides a useful approach as different uncertainties can~~
746 ~~be manipulated. However, virtual ecology has its own set of limitations. The data~~
747 ~~degradation and sampling processes described here still represent a relatively ideal~~
748 ~~situation in that the timespan of the data is long relative to the driving processes, and the~~
749 ~~sub-sampling treatment is at regular depth intervals. The representation of mixing is also~~
750 ~~simple in that it is applied consistently down the core.~~ Thus, investigation of uncertainty
751 ~~through both empirical and virtual approaches is necessary to better understand the~~
752 ~~influence of process and observer error on the analysis of palaeoecological data.~~ Empirical
753 approaches could involve collecting high-quality data (e.g., well-dated sediment cores with
754 frequent sub-sampling resolution), ideally with replicate cores from the same location, to
755 use as a benchmark against which to assess analyses when sub-sampling resolution is
756 reduced (e.g., Liu et al., 2012). Experimental approaches might include laboratory and *in*
757 *situ* manipulation; for example, Payne and Gehrels (2010) monitor the movement of tephra
758 in the field and under controlled laboratory environments to understand the influence of
759 tephra taphonomy on tephrochronology. [Finally, combined Semi-empirical approaches](#)
760 [would combining](#) empirical-virtual [approaches methods](#) could be developed using
761 (virtually) modified empirical data; for example, applying a simulated mixing process to
762 empirical sediment core data (such as those from varved sediments that are subject to
763 minimal mixing) and assessing the subsequent analyses. Mann and Rutherford (2002)
764 demonstrate this approach by generating pseudoproxy data by subjecting instrumental
765 data to degradation, such as various noise processes and spatial sampling strategies, to
766 assess sea surface temperature reconstruction methods. They used simulation to create a
767 continuous sea surface temperature record from the patchier instrumental data before
768 applying the degradation processes. Although the implications of such methods are
769 different for ecological data, similar approaches could fruitfully be applied.

770 771 5 Conclusion

772 [Palaeoecological uncertainty can be considered at four levels: environmental processes,](#)
773 [field methods, laboratory methods, and quantitative analyses \(Table 1\).](#) The effects of

774 different sources of uncertainty are challenging to disentangle and quantify from empirical
775 studies alone, and virtual ecology provides a useful approach as different uncertainties can
776 be manipulated. However, virtual ecology has its own set of limitations. The data
777 degradation and sampling processes described here still represent a relatively ideal
778 situation in that the timespan of the data is long relative to the driving processes, and the
779 sub-sampling treatment is at regular depth intervals [for example](#). Thus, investigation of
780 uncertainty through both empirical and virtual approaches is necessary to better
781 understand the influence of process and observer error on the analysis of palaeoecological
782 data. ~~A better understanding of the proxy system models of different proxies (i.e., how~~
783 ~~different proxies record environmental signals in an archive~~ A better understanding of the
784 how different proxies record environmental signals in an archive (e.g., how closely coupled
785 a signal recorded by a proxy data is to the environmental change) and the uncertainties
786 around quantifying and analysing proxy data can bring us closer to understanding long-
787 term climate and ecosystem dynamics. Although we have assessed sources of uncertainty
788 on pseudoproxy data representing species communities, all proxies such as isotopic data,
789 or tree ring data, have their own idiosyncratic sources of uncertainty. Virtual Ecological
790 approaches can help towards assessing their uncertainties.

791

792

793

794 **Acknowledgements**

795 The authors are very grateful for the support of the Centre for eResearch, [University of](#)
796 [Auckland](#) particularly Nick Young and Noel Zeng for their digital research expertise. We
797 also acknowledge the New Zealand eScience Infrastructure (NeSI) for the use of their
798 computational resources, without which the scale of the analyses involved in this paper
799 would not have been possible. Callum Walley and Anthony Shaw of NeSI provided
800 instrumental support for high-performance computing. Finally, we extend our thanks to
801 members of the Perry Lab for helping shape the project.

802 **Data availability**

803 Functions used for Feature Analysis for Time-series, and other calculated metrics can be
804 found here: <https://doi.org/10.5281/zenodo.8052806>.

805

806 **Author contributions**

807 Quinn Asena contributed to the conceptualization, methodology, and writing the original
808 draft. George Perry contributed to the conceptualization, methodology, and reviewing and
809 editing the manuscript, providing supervision and statistical expertise. Janet Wilmshurst
810 provided palaeoecological expertise to the project, and reviewed and edited the
811 manuscript.

847

848 **Funding**

849 The author(s) disclosed receipt of the following financial support for the research,
850 authorship, and/or publication of this article: This work was conducted through the New
851 Zealand's Biological Heritage National Science Challenge, which is funded by the New
852 Zealand Ministry of Business, Innovation and Employment.

853

854 **References**

- 855 Adesanya Adeleye, M., Charles Andrew, S., Gallagher, R., van der Kaars, S., De Deckker, P.,
856 Hua, Q., Haberle, S.G., 2023. On the timing of megafaunal extinction and associated floristic
857 consequences in Australia through the lens of functional palaeoecology. *Quaternary Science*
858 *Reviews* 316, 108263. <https://doi.org/10.1016/j.quascirev.2023.108263>.
- 859 Asena, Quinn, George LW Perry, and Janet M Wilmshurst. 2024. "Is the Past Recoverable
860 from the Data? Pseudoproxy Modelling of Uncertainties in Palaeoecological Data." *The*
861 *Holocene*, May, 09596836241247304. <https://doi.org/10.1177/09596836241247304>.
- 862 Asena, Quinn, Nick Young, and Alex Pletzer. 2023. "UoA-eResearch/fisher: V1.0.0." Zenodo.
863 <https://doi.org/10.5281/ZENODO.8052806>.
- 864 Beck, Kristen K., Michael-Shawn Fletcher, Patricia S. Gadd, Henk Heijnis, Krystyna M.
865 Saunders, Gavin L. Simpson, and Atun Zawadzki. 2018. "Variance and Rate-of-Change as
866 Early Warning Signals for a Critical Transition in an Aquatic Ecosystem State: A Test Case
867 from Tasmania, Australia." *Journal of Geophysical Research: Biogeosciences* 123 (2): 495–
868 508. <https://doi.org/10.1002/2017JG004135>.
- 869 Benito, Blas M., Graciela Gil-Romera, and H. John B. Birks. 2020. "Ecological Memory at
870 Millennial Time-Scales: The Importance of Data Constraints, Species Longevity and Niche
871 Features." *Ecography* 43 (1): 1–10. <https://doi.org/10.1111/ecog.04772>.
- 872 Birks, John B. H., André F. Lotter, Steve Juggins, and John P. Smol. 2012. *Tracking*
873 *Environmental Change Using Lake Sediments: Data Handling and Numerical Techniques*.
874 Springer Science & Business Media.
- 875 Blaauw, Maarten. 2012. "Out of Tune: The Dangers of Aligning Proxy Archives." *Quaternary*
876 *Science Reviews*, The INTegration of Ice Core, Marine and TERrestrial Records of the Last
877 Termination (INTIMATE) 60,000 to 8000 BP, 36 (March): 38–49.
878 <https://doi.org/10.1016/j.quascirev.2010.11.012>.
- 879 Blaauw, Maarten, K. D. Bennett, and J. Andrés Christen. 2010. "Random Walk Simulations of
880 Fossil Proxy Data." *The Holocene* 20 (4): 645–49.
881 <https://doi.org/10.1177/0959683609355180>.

Formatt

Formatt

- 882 Blaauw, Maarten, J. Andrés Christen, and Marco Antonio Aquino-López. 2020. "A Review of
883 Statistics in Palaeoenvironmental Research." *Journal of Agricultural, Biological and*
884 *Environmental Statistics* 25 (1): 17–31. <https://doi.org/10.1007/s13253-019-00374-2>.
- 885 Bothe, Oliver, Sebastian Wagner, and Eduardo Zorita. 2019. "Simple Noise Estimates and
886 Pseudoproxies for the Last 21000 Years." *Earth System Science Data* 11 (3): 1129–52.
887 <https://doi.org/10.5194/essd-11-1129-2019>.
- 888 Cabezas, Heriberto, and Brian D. Fath. 2002. "Towards a Theory of Sustainable Systems."
889 *Fluid Phase Equilibria*, Proceedings of the Ninth International Conference on Properties and
890 Phase Equilibria for Product and Process Design, 194–197 (March): 3–14.
891 [https://doi.org/10.1016/S0378-3812\(01\)00677-X](https://doi.org/10.1016/S0378-3812(01)00677-X).
- 892 Christiansen, Bo, T. Schmith, and P. Thejll. 2009. "A Surrogate Ensemble Study of Climate
893 Reconstruction Methods: Stochasticity and Robustness." *Journal of Climate* 22 (4): 951–76.
894 <https://doi.org/10.1175/2008JCLI2301.1>.
- 895 Clements, Christopher F., and Arpat Ozgul. 2016. "Including Trait-Based Early Warning
896 Signals Helps Predict Population Collapse." *Nature Communications* 7 (1).
897 <https://doi.org/10.1038/ncomms10984>.
- 898 De'ath, Glenn. 1999. "Principal Curves: A New Technique for Indirect and Direct Gradient
899 Analysis." *Ecology* 80 (7): 2237–53. [https://doi.org/10.1890/0012-9658\(1999\)080\[2237:PCANTF\]2.0.CO;2](https://doi.org/10.1890/0012-9658(1999)080[2237:PCANTF]2.0.CO;2).
- 901 Dormann, Carsten F., Jana M. McPherson, Miguel B. Araújo, Roger Bivand, Janine Bolliger,
902 Gudrun Carl, Richard G. Davies, et al. 2007. "Methods to Account for Spatial Autocorrelation
903 in the Analysis of Species Distributional Data: A Review." *Ecography* 30 (5): 609–28.
904 <https://doi.org/10.1111/j.2007.0906-7590.05171.x>.
- 905 Eason, Tarsha, and Heriberto Cabezas. 2012. "Evaluating the Sustainability of a Regional
906 System Using Fisher Information in the San Luis Basin, Colorado." *Journal of Environmental*
907 *Management* 94 (1): 41–49. <https://doi.org/10.1016/j.jenvman.2011.08.003>.
- 908 Eason, Tarsha, Ahjond S. Garmestani, Craig A. Stow, Carmen Rojo, Miguel Alvarez-Cobelas,
909 and Heriberto Cabezas. 2016. "Managing for Resilience: An Information Theory-Based
910 Approach to Assessing Ecosystems." *Journal of Applied Ecology* 53 (3): 656–65.
911 <https://doi.org/10.1111/1365-2664.12597>.
- 912 Evans, M. N., S. E. Tolwinski-Ward, D. M. Thompson, and K. J. Anchukaitis. 2013.
913 "Applications of Proxy System Modeling in High Resolution Paleoclimatology." *Quaternary*
914 *Science Reviews* 76 (September): 16–28. <https://doi.org/10.1016/j.quascirev.2013.05.024>.
- 915 Fisher, R. A. 1922. "On the Mathematical Foundations of Theoretical Statistics."
916 *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a*
917 *Mathematical or Physical Character* 222: 309–68.
- 918 Goring, Simon, Terri Lacourse, Marlow G. Pellatt, and Rolf W. Mathewes. 2013. "Pollen
919 Assemblage Richness Does Not Reflect Regional Plant Species Richness: A Cautionary Tale."

- 920 Edited by Amy Austin. *Journal of Ecology* 101 (5): 1137–45. <https://doi.org/10.1111/1365->
921 2745.12135.
- 922 Hastie, Trevor, and Werner Stuetzle. 1989. “Principal Curves.” *Journal of the American*
923 *Statistical Association* 84 (406): 502–16.
924 <https://doi.org/10.1080/01621459.1989.10478797>.
- 925 Jackson, Stephen T. 2007. “Looking Forward from the Past: History, Ecology, and
926 Conservation.” *Frontiers in Ecology and the Environment* 5 (9): 455–55.
927 [https://doi.org/10.1890/1540-9295\(2007\)5\[455:LFFTPH\]2.0.CO;2](https://doi.org/10.1890/1540-9295(2007)5[455:LFFTPH]2.0.CO;2).
- 928 Karunanithi, Arunprakash T., Heriberto Cabezas, B. Roy Frieden, and Christopher W.
929 Pawlowski. 2008. “Detection and Assessment of Ecosystem Regime Shifts from Fisher
930 Information.” *Ecology and Society* 13 (1).
- 931 Killick, Rebecca, and Idris A. Eckley. 2014. “Changepoint: An r Package for Changepoint
932 Analysis.” *Journal of Statistical Software* 58 (1): 1–19.
933 <https://doi.org/10.18637/jss.v058.i03>.
- 934 Kim, Dae-Won, Pavlos Protopapas, Yong-Ik Byun, Charles Alcock, Roni Khardon, and
935 Markos Trichas. 2011. “Quasi-Stellar Object Selection Algorithm Using Time Variability and
936 Machine Learning: Selection of 1620 Quasi-Stellar Object Candidates from Macho Large
937 Magellanic Cloud Database.” *The Astrophysical Journal* 735 (2): 68.
938 <https://doi.org/10.1088/0004-637X/735/2/68>.
- 939 Kosnik, Matthew A., and Michał Kowalewski. 2016. “Understanding Modern Extinctions in
940 Marine Ecosystems: The Role of Palaeoecological Data.” *Biology Letters* 12 (4).
941 <https://doi.org/10.1098/rsbl.2015.0951>.
- 942 Liu, Yao, Simon Brewer, Robert K. Booth, Thomas A. Minckley, and Stephen T. Jackson.
943 2012. “Temporal Density of Pollen Sampling Affects Age Determination of the Mid-
944 Holocene Hemlock (*Tsuga*) Decline.” *Quaternary Science Reviews* 45 (June): 54–59.
945 <https://doi.org/10.1016/j.quascirev.2012.05.001>.
- 946 Mann, Michael E., and Scott Rutherford. 2002. “Climate Reconstruction Using
947 ‘Pseudoproxies’.” *Geophysical Research Letters* 29 (10): 139-1-139-4.
948 <https://doi.org/10.1029/2001GL014554>.
- 949 Mayer, Audrey L., Christopher W. Pawlowski, and Heriberto Cabezas. 2006. “Fisher
950 Information and Dynamic Regime Changes in Ecological Systems.” *Ecological Modelling*,
951 Selected Papers from the Third Conference of the International Society for Ecological
952 Informatics (ISEI), August 26–30, 2002, Grottaferrata, Rome, Italy, 195 (1): 72–82.
953 <https://doi.org/10.1016/j.ecolmodel.2005.11.011>.
- 954 Nun, Isadora, Pavlos Protopapas, Brandon Sim, Ming Zhu, Rahul Dave, Nicolas Castro, and
955 Karim Pichara. 2015. “FATS: Feature Analysis for Time Series.” *arXiv:1506.00010*, August.
956 <https://arxiv.org/abs/1506.00010>.

- 957 Parnell, A. C., J. Haslett, J. R. M. Allen, C. E. Buck, and B. Huntley. 2008. "A Flexible Approach
958 to Assessing Synchronicity of Past Events Using Bayesian Reconstructions of Sedimentation
959 History." *Quaternary Science Reviews* 27 (19): 1872–85.
960 <https://doi.org/10.1016/j.quascirev.2008.07.009>.
- 961 Payne, Richard, and Maria Gehrels. 2010. "The Formation of Tephra Layers in Peatlands: An
962 Experimental Approach." *CATENA* 81 (1): 12–23.
963 <https://doi.org/10.1016/j.catena.2009.12.001>.
- 964 Peck, Steven L. 2004. "Simulation as Experiment: A Philosophical Reassessment for
965 Biological Modeling." *Trends in Ecology & Evolution* 19 (10): 530–34.
966 <https://doi.org/10.1016/j.tree.2004.07.019>.
- 967 Plumpton, Heather, Bronwen Whitney, and Francis Mayle. 2019. "Ecosystem Turnover in
968 Palaeoecological Records: The Sensitivity of Pollen and Phytolith Proxies to Detecting
969 Vegetation Change in Southwestern Amazonia." *The Holocene* 29 (1): 1720–30.
970 <https://doi.org/10.1177/0959683619862021>.
- 971 R Core Team. 2020. "R: A Language and Environment for Statistical Computing." Vienna,
972 Austria.
- 973 Reitalu, Triin, Pille Gerhold, Anneli Poska, Meelis Pärtel, Vivika Väli, and Siim Veski. 2015.
974 "Novel Insights into Post-Glacial Vegetation Change: Functional and Phylogenetic Diversity
975 in Pollen Records." Edited by Otto Wildi. *Journal of Vegetation Science* 26 (5): 911–22.
976 <https://doi.org/10.1111/jvs.12300>.
- 977 Richards, Joseph W., Dan L. Starr, Nathaniel R. Butler, Joshua S. Bloom, John M. Brewer,
978 Arien Crellin-Quick, Justin Higgins, Rachel Kennedy, and Maxime Rischard. 2011. "On
979 Machine-Learned Classification of Variable Stars with Sparse and Noisy Time-Series Data."
980 *The Astrophysical Journal* 733 (1): 10. <https://doi.org/10.1088/0004-637X/733/1/10>.
- 981 Roberts, Caleb P., Dirac Twidwell, Jessica L. Burnett, Victoria M. Donovan, Carissa L.
982 Wonkka, Christine L. Bielski, Ahjond S. Garmestani, et al. 2018. "Early Warnings for State
983 Transitions." *Rangeland Ecology & Management* 71 (6): 659–70.
984 <https://doi.org/10.1016/j.rama.2018.04.012>.
- 985 Simpson, G. L., and J. Oksanen. 2020. *Analogue Matching and Modern Analogue Technique*
986 *Transfer Function Models. Version 0.17-4*.
- 987 Simpson, Gavin L., and H. John B. Birks. 2012. "Statistical Learning in Palaeolimnology." In
988 *Tracking Environmental Change Using Lake Sediments: Data Handling and Numerical*
989 *Techniques*, edited by H. John B. Birks, André F. Lotter, Steve Juggins, and John P. Smol, 249–
990 327. Developments in Paleoenvironmental Research. Dordrecht: Springer Netherlands.
991 https://doi.org/10.1007/978-94-007-2745-8_9.
- 992 Smerdon, Jason E. 2012. "Climate Models as a Test Bed for Climate Reconstruction
993 Methods: Pseudoproxy Experiments." *WIREs Climate Change* 3 (1): 63–77.
994 <https://doi.org/10.1002/wcc.149>.

- 995 Sokolovsky, K. V., P. Gavras, A. Karampelas, S. V. Antipin, I. Bellas-Velidis, P. Benni, A. Z.
996 Bonanos, et al. 2017. "Comparative Performance of Selected Variability Detection
997 Techniques in Photometric Time Series." *Monthly Notices of the Royal Astronomical Society*
998 464 (1): 274–92. <https://doi.org/10.1093/mnras/stw2262>.
- 999 Spanbauer, Trisha L., Craig R. Allen, David G. Angeler, Tarsha Eason, Sherilyn C. Fritz,
1000 Ahjond S. Garmestani, Kirsty L. Nash, and Jeffery R. Stone. 2014. "Prolonged Instability
1001 Prior to a Regime Shift." Edited by John A. D. Aston. *PLoS ONE* 9 (10).
1002 <https://doi.org/10.1371/journal.pone.0108936>.
- 1003 Spanbauer, Trisha L., Craig R. Allen, David G. Angeler, Tarsha Eason, Sherilyn C. Fritz,
1004 Ahjond S. Garmestani, Kirsty L. Nash, Jeffery R. Stone, Craig A. Stow, and Shana M.
1005 Sundstrom. 2016. "Body Size Distributions Signal a Regime Shift in a Lake Ecosystem."
1006 *Proceedings of the Royal Society B: Biological Sciences* 283 (1833).
1007 <https://doi.org/10.1098/rspb.2016.0249>.
- 1008 Taranu, Zofia E., Stephen R. Carpenter, Victor Frossard, Jean-Philippe Jenny, Zoë Thomas,
1009 Jesse C. Vermaire, and Marie-Elodie Perga. 2018. "Can We Detect Ecosystem Critical
1010 Transitions and Signals of Changing Resilience from Paleo-Ecological Records?" *Ecosphere*
1011 9 (10). <https://doi.org/10.1002/ecs2.2438>.
- 1012 Telford, R.J., Heegaard, E., Birks, H.J.B., 2004. All Age-Depth Models Are Wrong: But How
1013 Badly? *Quaternary Science Reviews* 23, 1–5.
1014 <https://doi.org/10.1016/J.QUASCIREV.2003.11.003>
- 1015 Williams, John W., Jessica L. Blois, and Bryan N. Shuman. 2011. "Extrinsic and Intrinsic
1016 Forcing of Abrupt Ecological Change: Case Studies from the Late Quaternary." *Journal of*
1017 *Ecology* 99 (3): 664–77. <https://doi.org/10.1111/j.1365-2745.2011.01810.x>.
- 1018 [Wilmshurst, Janet M., Matt S. McGlone, and Dan J. Charman. "Holocene Vegetation and](#)
1019 [Climate Change in Southern New Zealand: Linkages between Forest Composition and](#)
1020 [Quantitative Surface Moisture Reconstructions from an Ombrogenous Bog." *Journal of*](#)
1021 [Quaternary Science](#) 17, no. 7 (2002): 653–66. <https://doi.org/10.1002/jqs.689>.
- 1022 Zurell, Damaris, Uta Berger, Juliano S. Cabral, Florian Jeltsch, Christine N. Meynard, Tamara
1023 Münkemüller, Nana Nehrbass, et al. 2010. "The Virtual Ecologist Approach: Simulating Data
1024 and Observers." *Oikos* 119 (4): 622–35. [https://doi.org/10.1111/j.1600-](https://doi.org/10.1111/j.1600-0706.2009.18284.x)
1025 [0706.2009.18284.x](https://doi.org/10.1111/j.1600-0706.2009.18284.x).