#### Fossilization potential of marine assemblages and environments 1

# 2

- Jack O. Shaw<sup>1\*</sup>, Derek E. G. Briggs<sup>1,2</sup>, and Pincelli M. Hull<sup>1,2</sup> <sup>1</sup>Department of Earth and Planetary Sciences, Yale University, New Haven, Connecticut 3 4 06511, USA
- <sup>2</sup>Peabody Museum of Natural History, Yale University, New Haven, Connecticut 06520, 5 USA
- 6
- 7

#### Table of Contents 8

9	SUPPLEMENTARY INFORMATION
10	DATA PROCESSING
11	Pelagic and seamount assemblages
12	TAXON- AND GEOGRAPHY-SPECIFIC ANALYSES
13	Sampling biases4
14	TAXON DURATIONS AND FOSSILIZATION POTENTIAL
15	INFERRED VERSUS STRATIGRAPHIC TAXON DURATIONS
16	Predictive modeling
17	IMPACT OF KONSERVAT-LAGERSTÄTTEN ON INFERRED FOSSILIZATION POTENTIAL OF TAXA
18	SUPPLEMENTARY FIGURES10
19	SUPPLEMENTARY DATA DESCRIPTIONS
20	REFERENCES
21	

### 23 SUPPLEMENTARY INFORMATION

24

## 25 Data processing

We downloaded all metazoan occurrences from the OBIS on 21<sup>st</sup> January 2020 (Supplementary 26 27 Data DR1) and retained only those species names recognized in the World Register of Marine 28 Species (WoRMS) (Horton et al., 2020). We removed occurrences (i) lacking data on genus, 29 latitude, longitude, depth, or duration of sampling, (ii) from non-marine settings, or (iii) with 30 inadequate error control (e.g., citizen science data). Occurrences were grouped into assemblages 31 if they were obtained within 0.1° latitude and 0.1° longitude of each other, in the same 10 m 32 depth bin, during the same year, and from the same dataset. We retained only assemblages with 33 representatives of at least three phyla, four classes, five orders, six families, and seven genera in 34 order to avoid poorly sampled assemblages (Supplementary Data DR2). Alternative binning 35 procedures were tested and did not affect interpretations (Fig. S1). Data processing and analyses 36 were performed in R. All data are provided in the Data Repository (see pages 27-29).

Environmental assignments (shallow = continental shelf and above, deep = continental slope and below, coral reef, and seamount) were based on dataset descriptions and not inferences from faunal lists. We assigned substrate type (mud, sand, gravel, rock) to shallow assemblages (n = 7,545) based on data interpolated using K-nearest neighbor classification algorithms applied to the dbSEABED database (Jenkins, 2008) (Fig. S6). We did not assign substrates to other environments (e.g., deep water, seamounts, etc.) due to greater uncertainties introduced by the combined effects of lower data densities, interpolation method, and grid size.

We downloaded fossil occurrences from the PBDB on 17<sup>th</sup> January 2020 (Supplementary Data 44 45 DR3). PBDB occurrences could not be assigned to substrate types due to a lack of sufficient 46 PBDB information on substrate associations. PBDB occurrences were assigned to environments 47 comparable to those in OBIS for the 25,732 genera with accompanying environmental data 48 (Supplementary Data DR6). We used data on Konservat-Lagerstätten (i.e., exceptionally 49 preserved fossil deposits listed in Muscente et al., 2017) which include 20,987 fossil occurrences 50 of 4,398 genera only 208 of which have living representatives recorded in our subset of OBIS 51 occurrences (n = 626,509) in order to test their impact on fossilization potential.

#### 53 Pelagic and seamount assemblages

54 In the case of rarely preserved pelagic and seamount assemblages, low within-environment

55 fossilization potentials reflect the rarity of these habitats in the rock record. Neither setting is

56 recorded directly in the PBDB—we had to rely on the presence of 'seamount' and 'pelagic' in

- 57 the data descriptions associated with PBDB occurrences. We identified seven geologic
- formations associated with seamounts in the PBDB, containing 483 occurrences and 213 genera.
- 59 Seventy-eight percent of the OBIS genera exclusive to seamounts (253 of 5,635) are represented

60 in the fossil record. Fifty-eight geologic formations in the PBDB are associated with dominantly

61 pelagic taxa, containing 3,437 occurrences and 1,359 genera. Nine percent of the OBIS genera

62 exclusive to pelagic assemblages (274 of 5,635) are represented in the fossil record. However,

63 OBIS includes only 6 pelagic datasets (Fig. S8) and they include predominantly fish and/or

64 zooplankton records and omit the abundant gelatinous members of the marine plankton. Thus the

65 true taxon fossilization potential of pelagic environments is likely closer to their within-

66 environment fossil potential of zero than that calculated here (Fig. 3, Fig. S3).

#### 67 Taxon- and geography-specific analyses

68 Given the large focus on invertebrates in the PBDB, we performed analyses of fossilization

- 69 potential excluding vertebrates from OBIS assemblages (Fig. S15). We found minimal changes
- 70 (<3%) to mean taxon and within-environment fossilization potential at the global level and when

71 data are parsed by substrate. There are also minimal changes (<5%) in shallow, deep, seamount,

72 and pelagic mean fossilization potential values. Coral reefs have higher invertebrate-only taxon

73 (from 44% to 54%) and within-environment (from 26% to 36%) mean fossilization potential

74 values (Fig. S4).

75 In order determine how the spatial distribution of sampling and taxa—for both OBIS and PBDB

76 data—impacted within-environment fossilization potential we analyzed a highly sampled area.

77 Few areas are both well sampled and include assemblages covering a range of environment

parameters, but we focused on an area with high PBDB and OBIS coverage, encompassing part

of Europe (bounded by coordinates: 58°N, 15°W; 58°N, 22°E; 35°N, 15°W; 35°N, 22°E). Our

80 sampling area contained areas of land and ocean in order to sample fossil and modern diversity.

81 The area contained 8,945 OBIS assemblages, of which 5512 included information allowing

82 within-environment fossilization potential to be calculated. All assemblages with environmental

83 information in the bounded area were shallow or deep water assemblages. Mean taxon 84 fossilization potential was 34% and mean within-environment fossilization potential 30% 85 (compared to 38% and 29%, respectively, for global data). Mean substrate taxon fossilization 86 potential was lowest in rock and gravel-based assemblages (24% and 20%, respectively) and 87 greatest in sand and mud-based assemblages (33% and 34%, respectively), conforming to the 88 same trend of higher fossilization potential in finer grained substrates evident in the global data. 89 Shallow water mean taxon fossilization potential was 33% and mean within-environment 90 fossilization potential was 31% (compared to 34% and 32%, respectively, for global data). Deep 91 water mean taxon fossilization potential was 27% and mean within-environment fossilization 92 potential was 14% (compared to 34% and 15%, respectively, for global data). The greatest 93 difference between the mean fossilization potential of global and European assemblages is the 94 taxon fossilization potential of deep water assemblages, which is 7% lower in the European 95 sample than the global sample (Fig. S9). This difference is unsurprising, given the large variation 96 in deep water taxon fossilization potential. The consistencies between the European subset and 97 the global data indicate that much of the trend in fossilization potential is robust to geographic 98 sampling biases.

#### 99 Sampling biases

OBIS assemblages and PBDB taxa are dominated by shallow-water representatives, and these
 OBIS assemblages generally show the highest fossilization potential values. To test the
 relationship between sampling and fossilization potential we performed a number of
 subsampling analyses. It is important to note that the fossilization potential metrics used here are
 inherently reliant on taxonomic diversity, such that subsampling would remove a key point of
 our study: the PBDB (and the fossil record) is biased against the fossilization of certain habitats.

We tested biases in sampling using the three best sampled groups in terms of number of OBIS assemblages and number of PBDB genera: shallow, coral reef, and deep water environments. We excluded seamount and pelagic communities because subsampling of such small sample sizes in terms of both OBIS assemblages and PBDB taxa—would obscure any trends. The smallest number of assemblages by environment is ~170 (seamounts) and the smallest number of fossil taxa by environment is ~200 (seamounts). 112 We subsampled to the smallest number of OBIS assemblages by environment (1012, coral reefs) 113 and subsampled to the smallest number of PBDB genera by environment (5927, deep water) to 114 calculate mean fossilization potential values (Fig. S13). We repeated this process 1000 times to 115 generate 1000 mean taxon and within-environment fossilization potential values for each 116 environment. Subsampling data in this way showed that (1) mean taxon fossilization potential 117 values generated using all assemblages (e.g., pictured in Fig. 3) are within the interquartile 118 ranges of mean values of subsampled data (Fig. S13) when parsed by environment, and (2) mean 119 within-environment fossilization potential values generated using all assemblages are greater 120 than mean values of subsampled data due to artificially reduced fossil diversity in the latter, and 121 (3) the larger the reduction in PBDB sample size, the greater the difference between all data 122 mean within-environment fossilization potential values and subsampled data mean values. The 123 within-environment fossilization potential calculation no longer serves its initial purpose of 124 illustrating contrast between PBDB environments, as fossil diversity is maintained for the least 125 sampled environment group but reduced for the other groups.

126 We also varied the number of PBDB taxa sampled by environment (as opposed to sampling to

127 the smallest number of PBDB genera by environment, as described above) in order to consider

128 how biases in fossil sampling impact within-environment fossilization potential values. This

129 showed that, for a given level of fossil sampling, coral reef taxa are proportionally better

represented in the PBDB compared to shallow and deep water fossil taxa (Fig. S14).

131 Fossilization potential values at the PBDB sample size of 5927 (Fig. S14, x-axis) are the same as

those depicted in within-environment boxplots in Fig. S13.

#### 133 **Taxon durations and fossilization potential**

134 The duration of a taxon impacts the chances of it entering fossil record. A long-lived taxon is

more likely to enter the fossil record than a short-ranging taxon as it has a greater chance of

136 recovery. Taxon duration is thought to be determined primarily by geographic range (Jablonski

137 and Hunt, 2006). Additionally proposed biotic and abiotic determinants of duration include

habitat depth (Fortey, 1980; Jablonski and Bottjer, 1983), depth range, body size (Jablonski,

139 2008), life mode (Crampton et al., 2010), niche breadth (e.g., number of environments occupied:

140 Kammer et al., 1997). Thus we expect taxon duration to vary between environments and to

141 contribute to assemblage-specific fossilization potential.

142 We considered the influence of genus duration on fossilization potential in two ways: (1)

143 comparing the durations of OBIS genera recorded in the PBDB between environments and

144 substrates; and (2) including the assemblage-specific means of these durations in a predictive

145 model of genus-level fossilization potential. The mean genus duration of an assemblage only

146 included taxa recorded in the PBDB and excluded taxa without a fossil record.

147 To generate genus durations we utilized the Adaptive Beta method, which is designed to estimate

the true stratigraphic range of a taxon based on the temporal distribution of its fossil occurrences 149

(Wang et al., 2016). The code allowed only 150-200 occurrences per taxon to be included. For

150 taxa with more than 150 occurrences in the PBDB we subsampled 150 occurrences (always

151 including the oldest and youngest) 1000 times to generate 1000 estimated origination ages, and

152 used the mean of these ages as the inferred origination age (Data Repository). For taxa with

153 fewer than six occurrences we used the oldest occurrence as the origination age, and for taxa

154 with fewer than 150 occurrences we applied the method without subsampling.

155 We assembled a subset all of the OBIS genera recorded in communities with environmental

156 information (n = 9,669). We plotted these data in two forms: (1) boxplots of genus durations by

157 occurrence (Fig. S10A; i.e., genera recorded in multiple OBIS environments were counted

158 multiple times), and (2) boxplots of the assemblage-specific average genus durations (Fig.

159 S10B).

148

160 We found that genus durations were similar across substrate types (Fig. S10A). Genus durations

161 were also greatest in shallow water assemblages and lowest in seamount assemblages. When the

162 genus durations of taxa within assemblages were averaged we found that shallow water

163 assemblages have the highest mean values, whereas pelagic and seamount assemblages have the

164 lowest values. Mean assemblage-specific genus durations were greater in finer substrates than in

165 coarser substrates.

166 Using linear regressions, we found a significant—albeit weak—negative correlation between

167 assemblage-specific mean genus duration and taxon fossilization potential (slope = -0.04, p-

168 value << 0.05, adjusted R-squared = 0.04). We also found a significant, weak, positive

169 correlation between mean genus duration and within-environment fossilization potential (slope =

170 0.02, p-value << 0.05, adjusted R-squared = 0.01). When parsed by environment or substrate we found variable weakly positive and weakly negative correlations between genus duration andfossilization potential.

#### 173 Inferred versus stratigraphic taxon durations

174 The difference between inferred taxon duration (calculated using the Adaptive Beta method) and 175 stratigraphic taxon duration (calculated using the oldest known occurrence in the PBDB) is 176 determined by the relative density of fossil occurrences—the greater the density of occurrences, 177 the smaller the difference between inferred and stratigraphic taxon durations. Thus, any two 178 environments might differ in average inferred taxon duration because either (1) true taxon 179 duration, and thus inferred taxon duration, differs, or (2) true taxon duration is the same but 180 preservation differs, and thus the estimates of inferred taxon duration differ. In the first case we 181 would expect the sites with the best-preserved fossil assemblages to show the largest differences 182 between inferred and stratigraphic durations. In the second case we would expect the best-183 preserved sites to show the smallest differences between inferred and stratigraphic durations.

We find that the difference between inferred genus duration and stratigraphic duration are
greatest in shallow environments, followed by coral reefs and deep water environments,
seamounts, and pelagic environments (Fig. S16). In other words, differences are greatest in the
best-preserved environments, suggesting that we are identifying true differences in taxon
durations between environments, rather than differences generated by sampling artifacts.
Therefore we used inferred taxon duration as a proxy for true taxon duration in our analyses.

## 190 **Predictive modeling**

191 We used conditional inference classification trees (CTree) built in the R package *partykit* 

192 (Hothorn and Zeileis, 2015) to assess the relative importance of environmental differences in

193 preservation potential (variables: environment, water depth, substrate, realm, alpha diversity),

194 taxon longevity (variables: average genus duration), and sampling biases (variables: latitude,

- 195 longitude, alpha diversity) on the fossilization potential of marine assemblages (Data
- 196 Repository). CTree is a type of Classification and Regression Tree (CART) analysis that
- 197 iteratively tests subsets of data until the null hypothesis of interdependence among variables
- 198 cannot be rejected. Unlike some other classification tree models, it accounts for covariation
- among predictor variables for unbiased selection and identification of variable importance.

200 Additionally, the model was chosen because it handles overfitting, non-linear and non-

- 201 parametric data, and missing information. Models were built using only communities with
- 202 environmental data (i.e., shallow, coral reef, deep water, seamount, and pelagic communities; n =
- 9,738). We ran CTree analysis 1000 times for both types of fossilization potential. Relative
- 204 variable importance scores, R-squared (R2), and root mean square error (RMSE) values were
- then averaged across the 1000 runs. Each run utilized a random sample of 70% of all
- assemblages. Variable significance was assessed at an alpha of 0.05 using Bonferroni adjusted
- 207 Monte-Carlo p-values (10 permutations per analysis).
- 208 The dominance of spatial variables (longitude and latitude) in predictive models of fossilization
- 209 potential is indicative of sampling bias. This is most obvious in analyses of pelagic community
- 210 fossilization potential where different datasets focus on different taxa with very different
- 211 fossilization potentials (Fig. S8). However, this sampling bias does not necessarily erroneously
- skew fossilization potential calculations. Most datasets are inherently spatially and
- 213 environmentally restricted (i.e., one research group will focus on a small spatial patch, and not
- sample the globe), and these attributes covary (e.g., the presence of coral reefs around Australia).

#### 215 Impact of Konservat-Lagerstätten on inferred fossilization potential of taxa

- Taxa in Konservat-Lagerstätten add significant data to diversity and disparity estimates and our
  understanding of macroevolution particularly in the early Paleozoic (Briggs et al., 1994; Hou et
  al., 2017). Thus Konservat-Lagerstätten might be expected to have a dramatic effect on
- 219 inferences of fossilization potential by providing rare glimpses of ancient assemblages. However,
- just 2.2% of the 58,870 genera in the PBDB are known only from Konservat-Lagerstätten.
- 221 Furthermore, exceptionally preserved assemblages are much rarer in the Cenozoic (7 Konservat-
- Lagerstätten), when 75% (n=1,921) of OBIS genera with PBDB occurrences are found first, than
- in the Paleozoic (142 Konservat-Lagerstätten). Only 6 genera (0.05%, n=9,806) (7 families and 3
- orders) in OBIS are recorded only in Konservat-Lagerstätten deposits and not in other PBDB
- deposits. When they are omitted from calculations of fossilization potential (i.e., limiting the
- study to 'typical' preservation conditions), there is no difference in our findings (Fig. S1). Thus
- 227 our calculations of fossilization potential are robust to the inclusion (or exclusion) of Konservat-
- 228 Lagerstätten from the PBDB the shelly fossil record provides the primary archive of potential

- assemblage composition. Particularly in the Cenozoic, there are few Konservat-Lagerstätten to
- 230 fill known gaps in these biased archives.

# 232 SUPPLEMENTARY FIGURES





236

237 Supplementary Figure S1: Comparisons of taxon and within-environment fossilization 238 potential distributions calculated using alternative binning procedures. Unless otherwise noted distributions include the default criteria used in the main text: all PBDB data and OBIS 239 240 assemblages grouped by 10 m depth bins, 0.1° latitude-longitude bins, sampling year, and a 241 minimum diversity of three phyla, four classes, five orders, six families, and seven genera. 242 Assemblages grouped with (A) 20 m depth bins and 1° latitude-longitude bins, but not by dataset 243 ID); (B) year excluded; (C) 20 m depth bins; (D) 100 m depth bins; (E) 0.01° latitude-longitude 244 bins; (F) 1° latitude-longitude bins; (G) no minimum diversity criteria; (H) Lagerstätten excluded; (I) assemblages grouped using the default criteria (Genus panel: Fig. 1, main text). 245





249 **Supplementary Figure S2:** Taxon fossilization potential distributions for shallow-water

assemblages with substrate information (n = 7,545). Taxon fossilization potential calculated at

ordinal, family, and generic ranks. Histogram bin widths = 5%, dashed lines indicate mean.



253

Supplementary Figure S3: Taxon (red) and within-environment (blue) fossilization potential distributions of assemblages with environmental information (n = 9,669). Fossilization potential calculated at order, family, and genus ranks. Histogram bin widths = 5%, dashed lines indicate mean.



260Supplementary Figure S4: Percentage of taxa (orders, families, or genera) in assemblages261belonging to the phylum Chordata, plotted against taxon and within-environment fossilization262potential. Datapoints indicate assemblages (n = 9,669) and are colored by environment. Solid263lines indicate significant linear regressions (p-value < 0.05), dashed lines indicate insignificant</td>264regressions.



267 **Supplementary Figure S5:** Depth versus taxon (top) and within-environment (bottom)

- 268 fossilization potential for data from the BIOCEAN dataset (n = 549; detailed in main text). Line
- shows statistically significant linear regression with slope values indicated.



**Supplementary Figure S6:** Map of shallow water (<1000m) substrate type interpolated using

273 dbSeabed data and KNN algorithm.



- 276 Supplementary Figure S7: Bar chart displaying the percentages of OBIS genera with records in
- 277 the Paleobiology Database (PBDB). The total numbers of OBIS genera are indicated in
- 278 parentheses.
- 279
- 280







283 ID (n = 6) for all pelagic communities (n = 327).



Supplementary Figure S9: Taxon and within-environment fossilization potential distributions
for assemblages within a bounded area in Europe (58°N, 15°W; 58°N, 22°E; 35°N, 15°W; 35°N,
22°E). Histogram bin widths = 5%; dashed lines indicate mean.



	Lower whisker	Lower hinge	Median	Upper hinge	Upper whisker
Shallow	0.01	24.36	64.06	119.12	261.10
Coral reef	0.01	16.60	47.31	93.64	207.05
Deep	0.01	8.47	51.90	96.89	228.06
Pelagic	0.01	13.79	44.55	91.23	207.05
Seamount	0.01	1.36	41.04	85.90	207.05
Rock	0.01	18.45	64.56	120.03	270.48
Gravel	0.01	21.74	70.17	129.77	289.32
Sand	0.01	17.32	55.63	107.12	240.76
Mud	0.01	23.54	65.65	128.57	283.93

293 Supplementary Figure S10A: Boxplots (top) and boxplot statistics (bottom) of genus durations

294 (calculation method outlined in supplementary information section "Taxon durations and

295 fossilization potential") for OBIS taxa with corresponding environment and substrate

296 information. Taxa without fossil evidence are excluded. Outliers removed for clarity.



	Lower	Lower hinge	Median	Upper hinge	Upper
	whisker				whisker
Shallow	22.95	140.81	178.16	219.43	337.13
Coral reef	5.08	54.30	71.19	89.94	142.76
Deep	0.01	56.60	88.24	121.07	215.55
Pelagic	3.49	25.02	43.52	120.13	261.92
Seamount	16.78	36.74	57.06	89.51	148.55
Rock	0.01	115.03	156.48	315.54	404.18
Gravel	0.01	117.52	147.91	197.38	294.72
Sand	14.39	132.00	171.75	213.15	334.38
Mud	33.15	144.18	180.40	218.34	324.65

301 Supplementary Figure S10B Boxplots (top) and boxplot statistics (bottom) of mean genus

302 duration values for assemblages with environment (top left) and substrate (top right) information.

303 Calculations of assemblage mean values only include the durations of taxa with fossil

304 representatives. Outliers removed for clarity.



306

307 Supplementary Figure S11: Scaled variable importance scores for conditional inference

308 regression tree models predicting genus-level taxon and within-environment fossilization

309 potential. Importance rankings are based on the reduction in square error produced when the

310 variable is added. R-squared (R2) and root mean square error (RMSE) shown for each model.



313 Supplementary Figure S12: Assemblage-specific genus richness for assemblages with

environment (n = 9,669) and substrate (n = 7,545) data plotted against fossilization potential.



315

**Supplementary Figure S13:** Taxon and within-environment fossilization potential mean values for subsampled data (repeated 1000 times). Each mean represents comparisons utilizing random samples of PBDB taxa (5927 per environment, as per the smallest group, deep water) and random samples of OBIS assemblages (1012 per environment, as per the smallest group, coral reefs). Means of the raw data (e.g., Fig. 3) shown using black asterisks. Within-environment fossilization potential for shallow and coral reef environments are off the axes (indicated by circles) and have means of 32% and 26%, respectively.



325 Supplementary Figure S14: 95% confidence intervals of mean within-environment 326 fossilization potential values when calculated based on subsampled PBDB data. Mean within-327 environment fossilization potential values were generated by first subsampling the number of 328 PBDB taxa per environment (x-axis), then comparing the number of OBIS taxa in an assemblage 329 to the subsampled PBDB faunal list to generate the assemblage-specific fossilization potential 330 value, and finally calculating the mean of all assemblage-specific fossilization potentials by 331 environment. This process was repeated 1000 times at a number of pre-defined PBDB sample 332 sizes (100, 500, 1000, 2500, 5000, 7500, 10000, 20000) up to the maximum number of PBDB 333 taxa assigned to each environment.





336 Supplementary Figure S15: Invertebrate only analyses. (A) Taxon fossilization potential and

337 within-environment fossilization potential (mean taxon/within-environment fossilization

338 potential = 39/30%). (B) Taxon fossilization potential distributions for shallow water

assemblages with substrate information (mean fossilization potential values: rock = 22%, gravel

340 = 25%, sand = 36%, mud = 35%). (C) Taxon and within-environment fossilization potential

341 distributions for assemblages with environmental information (mean taxon/within-environment

- fossilization potential values: shallow = 33/30%, coral reef = 53/35%, deep = 35/14%, pelagic =
- 343 14/2%, seamount = 57/0%). Histogram bin widths = 5%; dashed lines indicate mean.



346 Supplementary Figure S16: Density plot of differences between assemblage averaged inferred-

347 and PBDB-durations parsed by environment.

### 348 SUPPLEMENTARY DATA DESCRIPTIONS

- 349 Supplementary data are available for download at the Harvard Dataverse Repository
- 350 (https://doi.org/10.7910/DVN/MMMTYZ). Supplementary Data labels and descriptions apply to
- data files listed in the Dataverse Repository.

Supplementary Data DR1: Raw list of all occurrences in OBIS downloaded on 21<sup>st</sup> January
 2020. Occurrences were culled (procedures in text) using R code (Sup. Data DR4A). Column

headers following Darwin Core format (Wieczorek et al., 2012).

355 **Supplementary Data DR2:** Raw list of OBIS datasets (downloaded 20<sup>th</sup> January 2020, Sup.

356 Data 2A) and list of OBIS assemblages (Sup. Data DR2B) meeting the culling procedures

described herein and applied using R Code (Sup. Data DR4A).

358 To generate assemblages, we removed occurrences from Sup. Data 1 (i) lacking data on 359 genus, latitude, longitude, depth, or duration of sampling, (ii) from non-marine settings, or (iii) 360 with inadequate error control (e.g., citizen science data). Occurrences were grouped into 361 assemblages if they were obtained within 0.1° latitude and 0.1° longitude of each other, in the 362 same 10 m depth bin, during the same year, and from the same dataset. We retained only 363 assemblages with representatives of at least three phyla, four classes, five orders, six families, 364 and seven genera in order to avoid poorly sampled assemblages. Alternative binning procedures 365 were tested and did not affect interpretations (Figs. S1).

366 Additional assemblage data (id habitat, id scope, id plankton, id sampling, id fishery, 367 and id environment) in Sup. Data DR2B was manually assigned to datasets using information 368 from dataset titles and abstracts in Sup. Data DR2A. Column headers in Sup. Data DR2B are as 369 follows: "dataset.par" = unique depth-, year-, location-, and dataset-specific assemblage 370 identifier; "dataset id" = an identifier for the set of data, which may be a global unique identifier 371 or an identifier specific to a collection or institution; "depth.slice" = 10m depth bin of 372 assemblage; "LongBin" =  $0.1^{\circ}$  longitudinal bin of assemblage; "LatBin" =  $0.1^{\circ}$  latitudinal bin of 373 assemblage; "id habitat" = marine realm of dataset, pelagic or benthic; "id scope" = taxonomic 374 scope of dataset; "id plankton" = plankton-only dataset; "id sampling" = sampling type of 375 assemblage; "id environment" = environment type of dataset; "title" = name of dataset; "abstract" = abstract of dataset. 376

377 Supplementary Data DR3: List of all occurrences in PBDB with select columns, downloaded

378 on 17<sup>th</sup> January 2020 (A). List of Lagerstätten names generated by manually searching for names

in PBDB ("formation") based on the list of Lagerstätten in Muscente et al. (2017) (B).

380 Data downloaded using the following URL:

381 http://paleobiodb.org/data1.2/occs/list.csv?datainfo&rowcount&base\_name=Animalia&show=fu
382 ll,genus,env

Joz II,genius,env

383 Supplementary Data DR4: R Code used for data cleaning (A), identifying assemblages and

calculating fossilization potential (B), analyzing fossilization potential (C), generating genus
 duration estimates (D; Wang et al., 2016), applying genus duration estimate calculations to OBIS

taxa (E), and predicting fossilization potential using conditional inference tree models (F).

Supplementary Data DR5: Seafloor substrate data (Jenkins, 2008; A), code used to create
predictive model (B), and predictive model (native R format, ".RDS") applied to OBIS
assemblages to assess substrate type (C).

Substrate data (Sup. Data DR5A) used for interpolation were supplied by Jenkins
(personal communications): "wd\_vals" = water depth; "lat" = latitude, to two decimal places;
"lon" = longitude, to two decimal places; "substrate" = dominant substrate type (rock, gravel, sand, and mud).

394 We interpolated substrate data based on the dbSeabed dataset (Jenkins, 2008) using a k-395 nearest neighbors (KNN) algorithm. The KNN algorithm is a non-parametric supervised machine 396 learning algorithm used for classification and regression. The KNN predictive model was built 397 using the R package KKNN (Schliep, Hechenbichler, and Lizee, 2016). We performed leave-398 one-out cross-validation for computationally efficient training of the predictive model. We tested 399 K-values ranging 1-20 and all possible kernel types ("rectangular", "triangular", "epanechnikov", 400 "gaussian", "rank", and "optimal"). For training, dbSeabed data was split so that 80% of the data 401 was used to train the model and 20% to test it. Training identified a K-value of 12 and a 402 triangular kernel as the best parameters, although kernel type did not alter resultant 403 classifications. Our best model had a k-value of 12 and a triangular kernel shape, meaning that

404 the 12-nearest neighbors within a triangle shaped space of the focal grid cell were used to infer405 the substrate type assigned to that grid cell.

406 Prior to analysis, we extended the dbSeabed data by 10 degrees of longitude on either 407 side (i.e., so that the map spanned -190 to 190 degrees). This allowed the predictive model to 408 utilize data at the edge of its conventional range (-180 to 180 degrees). If this correction was not 409 applied, the model would not know that data at -179 degrees longitude and 179 degrees longitude 410 are proximal. To extend the map, we simply duplicated data in the range of -180 to -170 degrees 411 and 170 to 180 degrees by adding and subtracting, respectively, 360 degrees (i.e., -179 became 412 181 degrees longitude).

413 Supplementary Data DR6: PBDB environments matched to broader categories for comparison
414 with OBIS designations.

415 Supplementary Data DR7: Tests of sampling biases for shallow water, coral reef, and deep
416 water assemblages (see supplementary information section "Sampling biases" for more
417 information).

### 418 **REFERENCES**

- Briggs, D.E.G., Erwin, D.H., and Collier, F.J., 1994, The fossils of Burgess Shale: Washington
  DC, Smithsonian Institution Press.
- 421 Crampton, J.S., Cooper, R.A., Beu, A.G., Foote, M., and Marshall, B.A., 2010, Biotic influences
- 422 on species duration: interactions between traits in marine molluscs: Paleobiology, v. 36, p.
- 423 204–223, doi:https://doi.org/10.1666/09010.1.
- Fortey, R.A., 1980, Generic longevity in lower ordovician trilobites: Relation to environment:
  Paleobiology, v. 6, p. 24–31, doi:https://doi.org/10.1017/S009483730001246X.
- 426 Horton, T. et al., 2020, World Register of Marine Species (WoRMS):,

427 https://www.marinespecies.org (accessed January 2020).

- Hothorn, T., and Zeileis, A., 2015, partykit : A Toolkit for Recursive Partytioning: Journal of
  Machine Learning Research, v. 16, p. 3905–3909.
- 430 Hou, X., Siveter, D.J., Siveter, D.J., Aldridge, R.J., Cong, P., Gabbott, S.E., Ma, X., Purnell,
- 431 M.A., and Williams, M., 2017, The Cambrian Fossils of Chengjiang, China: John Wiley &
- 432 Sons, doi:https://doi.org/10.1002/9781118896372.
- Jablonski, D., 2008, Species Selection: Theory and Data: Annual Review of Ecology, Evolution,
  and Systematics, v. 39, p. 501–524,
- 435 doi:https://doi.org/10.1146/annurev.ecolsys.39.110707.173510.
- 436 Jablonski, D., and Bottjer, D.J., 1983, Soft-bottom epifaunal suspension-feeding assemblages in
- the Late Cretaceous: implications for the evolution of benthic paleocommunities., *in* Biotic
  interactions in Recent and fossil benthic communities, p. 747–812.
- 439 Jablonski, D., and Hunt, G., 2006, Larval ecology, geographic range, and species survivorship in
- 440 cretaceous mollusks: Organismic versus species-level explanations: American Naturalist, v.
- 441 168, p. 556–564, doi:https://doi.org/10.1086/507994.
- 442 Jenkins, C., 2008, dbSEABED: an information processing system for marine substrates:,
- 443 https://instaar.colorado.edu/~jenkinsc/dbseabed/ (accessed April 2020).

444	Kammer, T.W., Baumiller, T.K., and Ausich, W.I., 1997, Species longevity as a function of
445	niche breadth: Evidence from fossil crinoids: Geology, v. 25, p. 219-222,
446	doi:https://doi.org/10.1130/0091-7613(1997)025<0219:SLAAFO>2.3.CO;2.
447	Muscente, A.D. et al., 2017, Exceptionally preserved fossil assemblages through geologic time
448	and space: Gondwana Research, v. 48, p. 164–188,
449	doi:https://doi.org/10.1016/J.GR.2017.04.020.
450	Schliep, K., Hechenbichler, K., and Lizee, A., 2016, kknn: Weighted k-Nearest Neighbors:,
451	https://cran.r-project.org/web/packages/kknn/index.html (accessed April 2020).
452	Wang, S.C., Everson, P.J., Zhou, H.J., Park, D., and Chudzicki, D.J., 2016, Adaptive credible
453	intervals on stratigraphic ranges when recovery potential is unknown: Paleobiology, v. 42,
454	p. 240-256, doi:https://doi.org/10.1017/pab.2015.37.
455	Wieczorek, J., Bloom, D., Guralnick, R., Blum, S., Döring, M., Giovanni, R., Robertson, T., and
456	Vieglais, D., 2012, Darwin core: An evolving community-developed biodiversity data

457 standard: PLoS ONE, doi:https://doi.org/10.1371/journal.pone.0029715.