

An open occurrence dataset for European subterranean spiders

Giuseppe Nicolosi^{1,2*}, Adrià Bellvert^{1,3}, Isabel R. Amorim⁴, Gergely Balázs^{5,6}, Anna Biró^{5,6}, Paulo AV Borges⁴, Traian Brad⁷, Tommaso Cancellario⁸, Pedro Cardoso⁹, Lorenzo Cresi¹, Luis Crespo⁴, Anđela Čukušić¹⁰, Iva Čupić¹⁰, Teo Delić¹¹, Sophie Front¹², Adrián García-Villarín¹³, Donard Geci¹⁴, Thomas Hesselberg¹⁵, Joanna Kocot-Zalewska¹⁶, Peter Kozel^{17,18}, Josiane Lips¹², Nuria Macías-Hernández¹⁹, Tone Novak¹⁷, Pedro Oromí¹⁹, Paolo Pantini²⁰, Kaloust Paragamian²¹, Savvas Paragkamian^{21,22}, Diego Patiño-Sauma¹⁹, Nikola Pischitta¹⁰, Žak Pušnar¹¹, Milan Rezac²³, Carles Ribera²⁴, Afonso Rodrigues⁹, Tin Rožman^{10, 25}, Ignac Sivec²⁶, Mojmir Štangelj²⁶, Karla Tolić¹⁰, Maja Zagmajster¹¹, Martina Pavlek^{27,10‡}, Stefano Mammola^{1,28,29‡*}

*Correspondence: Giuseppe Nicolosi (giuseppe.nicolosi@unict.it), Stefano Mammola (stefano.mammola@cnr.it)

‡Shared last author

¹Molecular Ecology Group (MEG), Water Research Institute (IRSA), National Research Council (CNR), Corso Tonolli 50, 28922 Verbania Pallanza, Italy

²Department of Biological, Geological and Environmental Sciences, University of Catania, Via A.Longo 19, 95125 Catania, Italy

³Zoological Institute and Museum, University of Greifswald, Loitzer Str. 26, Greifswald 17489, Germany

⁴University of Azores, CE3C—Centre for Ecology, Evolution and Environmental Changes, Azorean Biodiversity Group, CHANGE —Global Change and Sustainability Institute, School of Agricultural and Environmental Sciences, Rua Capitão João d'Ávila, Pico da Urze, 9700-042 Angra do Heroísmo, Azores, Portugal

⁵Department of Systematic Zoology and Ecology, Institute of Biology, ELTE Eötvös Loránd University, Budapest, Hungary

⁶HUN-REN-ELTE-MTM Integrative Ecology Research Group

⁷Institutul de Speologie „Emil Racoviță”, Str. Clinicilor Nr. 5-7, 400006 Cluj-Napoca, Romania

⁸Centre Balear de Biodiversitat, Universitat de les Illes Balears, Palma, Spain

⁹CE3C—Centre for Ecology, Evolution and Environmental Changes, CHANGE —Global Change and Sustainability Institute, Faculty of Sciences, University of Lisbon, Campo Grande, 1749-016 Lisboa, Portugal

¹⁰Croatian Biospeleological Society, Zagreb, Croatia

¹¹University of Ljubljana, Biotechnical Faculty, Department of Biology, SubBioLab, Jamnikarjeva 101, 1000 Ljubljana, Slovenia

¹²Biospeleology Study Group, French Federation of Speleology, France

- 40 ¹³Department of Animal Biology, Ecology, Parasitology, Edaphology and Agrochemistry,
41 University of Salamanca, Salamanca, Spain
- 42 ¹⁴Department of Biology, Faculty of Mathematics and Natural Sciences, University of
43 Prishtina, Mother Teresa Street p.n., 10000 Prishtinë, Republic of Kosovo
- 44 ¹⁵Department of Biology, University of Oxford. South Parks Road, Oxford OX1 3EL
- 45 ¹⁶Upper Silesian Museum, Department of Natural History, pl. Jana III Sobieskiego 2, 41-
46 902 Bytom, Poland
- 47 ¹⁷Department of Biology, Faculty of Natural Sciences and Mathematics, University of
48 Maribor, Maribor 2000, Slovenia
- 49 ¹⁸Karst Research Institute ZRC SAZU, Postojna 6230, Slovenia
- 50 ¹⁹Department of Animal Biology, Edaphology and Geology, University of La Laguna, La
51 Laguna, Tenerife, 38206 Canary Islands, Spain.
- 52 ²⁰Museo Civico di Scienze Naturali Enrico Caffi, Bergamo, Italy
- 53 ²¹Hellenic Institute of Speleological Research, Irakleio, Crete, Greece
- 54 ²²Department of Biology, University of Crete, Heraklion, Crete, Greece
- 55 ²³Czech Agrifood Research Center, Drnovská 507/73, 161 00 Prague, Czechia
- 56 ²⁴Departament de Biologia Evolutiva, Ecologia i Ciències Ambientals and Institut de
57 Recerca de la Biodiversitat, Universitat de Barcelona, Av.Diagonal 643, 08028 Barcelona,
58 Spain
- 59 ²⁵Croatian Natural History Museum, Zagreb, Croatia
- 60 ²⁶Slovenian Museum of Natural History, Prešernova, 20, 1000 Ljubljana, Slovenia
- 61 ²⁷Ruder Boskovic Institute, Zagreb, Croatia
- 62 ²⁸Finnish Museum of Natural History, University of Helsinki, Helsinki, Finland
- 63 ²⁹NBFC, National Biodiversity Future Center, Palermo, Italy

77 **Abstract**

78 Spiders are remarkably diverse in caves and other subterranean habitats, where they play
79 key ecological roles as generalist predators and strongly influence local food webs. They
80 have been instrumental as model organisms for testing various eco-evolutionary
81 hypotheses. Furthermore, strictly subterranean species exhibiting narrow ranges and high
82 endemism are particularly significant for conservation planning and vulnerability
83 assessments. Although high-quality data are essential for research on and conservation of
84 subterranean spiders, such information remains scarce, especially regarding distribution
85 patterns. To help fill this gap, we screened the literature, unpublished records, and open
86 datasets to compile georeferenced occurrences of subterranean spiders across caves and
87 other subterranean habitats throughout Europe. Based on these data—and to illustrate one
88 potential application of the compiled dataset—we present the first prediction of subterranean
89 spider richness patterns across Europe using stacked species distribution models. The
90 European Subterranean Spider Dataset (ESSD) comprises 31,224 records of 637
91 subterranean-dwelling spider species (including morphospecies under description),
92 covering a range of information including taxonomy, locality details (such as location name,
93 country, geographic coordinates, type of subterranean habitat), and reference information
94 for each record. All variables are coded using the Darwin Core Standard, ensuring
95 interoperability with the Global Biodiversity Information Facility (GBIF) and other biodiversity
96 databases. By enabling integration with trait and phylogenetic resources, the ESSD provides
97 a robust framework to investigate the drivers and processes shaping subterranean
98 biodiversity, assess vulnerability to environmental change and anthropogenic pressures,
99 and guide future sampling to progressively reduce geographic and taxonomic gaps through
100 open data sharing.

101 **Keywords:** Araneae, Darwin Core Standard, Hypogean, Open data, Species distribution
102 modelling (SDM)

103

104 **Introduction**

105 In recent years, there has been an explosion of biogeography and (macro)ecology studies
106 focused on uncovering the patterns and factors that shape biodiversity patterns at
107 increasingly larger spatial scales—continental to global (e.g., Labouyrie et al., 2023;
108 Martínez-Núñez et al., 2023, Sabatini et al., 2022). Building this understanding
109 fundamentally relies on high-quality data, especially distribution records. Over the past few
110 decades, online biodiversity databases have experienced substantial growth, largely due
111 to collaborative efforts that have enhanced data accessibility and sharing. This progress
112 has led to the formation of comprehensive databases covering a wide range of taxa,
113 including the Global Biodiversity Information Facility (GBIF) (GBIF, 2025), LifeWatch ERIC
114 (LifeWatch ERIC, 2025) and BioTIME (Dornelas et al., 2025). These repositories provide
115 extensive taxonomic and distributional information on thousands of taxa across various
116 ecosystem types and time scales. Despite these advances, large gaps and biases in
117 species' known geographic distributions, the so-called Wallacean shortfall, continue to limit
118 the completeness and reliability of macroecological and biogeographic inferences (e.g.,
119 Cardoso et al., 2011; Hortal et al., 2015; Hughes et al., 2021).

120 Primarily due to accessibility challenges (Ficetola et al., 2019; Mammola et al.,
121 2021a), the documentation of biodiversity in subterranean ecosystems (caves,
122 groundwaters, fissural systems, and the like) has historically progressed more slowly than
123 at the surface. However, in recent years, a steady accumulation of knowledge, combined
124 with the funding of specific projects focused on continental biodiversity inventories (e.g.,
125 PASCALIS, Biodiversa+ DarCo), has led many authors to compile this information and
126 publish it in public datasets with varying resolutions and scales. For example, we now
127 have the first global datasets on the distribution of cave-dwelling bats (Tanalgo et al.,
128 2022), cave fish (Bai et al., 2025), asellids (Saclier et al., 2024), and microwhip scorpions
129 (Mammola et al., 2021b). In Europe, Pascalis dataset (Deharveng et al. 2009) and
130 European Groundwater Crustaceans Dataset (Zagmajster et al. 2014) were used for early
131 analyses of continental patterns, but publication of the continental datasets of distribution
132 of subterranean organisms have started only in the last few years, like for bats (Fialas et
133 al. 2025), copepods (Cerasoli et al., 2025) and ostracods (Mori et al., 2025). However, a
134 similar large-scale dataset is still lacking for subterranean spiders.

135 Spiders (Arachnida: Araneae), with over 53,000 species currently described (World
136 Spider Catalog, 2025), and providing numerous essential ecosystem services (Cardoso et
137 al. 2025), are among the most widespread and generalist predators in terrestrial habitats
138 (Turnbull, 1973). Spiders are particularly diversified in caves and other subterranean voids,
139 where they play a key ecological role as predators and strongly structure local food webs.
140 Despite the growing interest in subterranean spiders (Mammola and Isaia, 2017), major
141 knowledge gaps remain, especially regarding their distribution. Limited expertise and lower
142 research interest in certain regions have delayed comprehensive data collection. However,
143 recent efforts are beginning to address these gaps, contributing essential data for
144 advancing our understanding of subterranean spider ecology and biodiversity. These
145 efforts include the publication of trait data for all the species in Europe (Mammola et al.,
146 2022; Patiño-Sauma et al., 2025) and high-resolution distribution data for selected caves
147 (Mammola et al., 2019a; Macías-Hernández et al., 2024) or regions (e.g., Western Alps;
148 Nicolosi et al., 2025; Azores; Crespo et al. 2025).

149 We present a novel dataset comprising occurrence records for all known species
150 and morphospecies of subterranean spiders from Europe, spanning a wide range of
151 ecological affinities to subterranean habitats, ranging from species still able to exploit
152 surface habitats to obligate subterranean dwellers. This dataset is the result of a
153 collaborative effort among multiple partners, including ecologists, conservationists,
154 taxonomists, and biogeographers from various European countries and outermost islands
155 of Europe (Azores, Madeira, Selvagens, and the Canary Islands). Their contributions have
156 significantly enriched the data availability, culminating in a comprehensive, multi-species
157 dataset designed to advance research and conservation efforts for these species. By
158 making these data public, we hope to promote collaborative research on subterranean
159 spider biodiversity, spatial patterns, drivers of distribution patterns, and quantitative
160 conservation efforts.

161

162 **Methods**

163

164 ***Target species and habitats***

165 We focused on subterranean spiders across continental Europe, including the
166 archipelagos of the Azores, Salvagens, Madeira, and Canary Islands. Following the
167 function-based classification of Earth's ecosystems (Keith et al., 2020, 2022), we focused
168 on ecosystems belonging to the 'Subterranean' (S) domain, which includes diverse
169 terrestrial subterranean systems: i) the 'Subterranean lithic' (S1) biome, namely various
170 type of caves (e.g., aerobic caves, lava tubes, volcanic pits) and other subterranean voids
171 of smaller sizes (e.g., fissure systems, deep scree strata, and the so-called *Milieu*
172 *Souterrain Superficial* [MSS; reviewed in Mammola et al. (2016)]); and ii) the
173 'Anthropogenic subterranean voids' (S2) biomes, namely all anthropogenic subterranean
174 voids with cave-like environmental conditions, including mines, underground bunkers,
175 blockhouses, tunnels, culverts, and cellars.

176 For the list of target species, we used the latest checklist of European subterranean
177 spiders by Patiño-Sauma et al. (2025), currently listing 637 species across 28 families. Of
178 these, 64 are species under description (hereinafter 'morphospecies'), identified by experts
179 as new taxonomic entities based on morphological and/or genetic information.

180

181 ***Data acquisition***

182 The spider dataset is a comprehensive compilation of diverse data sources, created
183 through the collaboration of 40 researchers across Europe. For each of the target species,
184 we compiled occurrence records based on different sources. First, we mined the primary
185 literature for reported localities and georeferenced records. Second, we included
186 unpublished records (e.g., data stored in personal and institutional collections) as provided
187 by each contributor in the authors list. Third, we mined records from the main accessible
188 databases on spiders in Europe, namely ArachnoMap (de Biurrun et al., 2022), Araneae.it
189 (Pantini & Isaia, 2019), the UK Spider and Harvestman Recording Scheme
190 (<https://srs.britishspiders.org.uk/>), Canary Islands Biodiversity Database (Gobierno de
191 Canarias, 2024), the Cave fauna of Greece database (Paragamian et al., 2025) and
192 regional database on subterranean species of the Western Balkans SubBioDB
193 (Zagmajster 2016). Lastly, we mined the GBIF database, which yielded a great number of
194 missing records, especially for the most widespread species. We downloaded the
195 occurrence records from GBIF using the Python package "biodumpy" v.0.1.6 (Cancellario
196 et al., 2025). Specifically, we employed the GBIF module, setting the parameter
197 `dataset_key` to "d7dddbf4-2cf0-4f39-9b2a-bb099caae36c" and `geometry` to "POLYGON((-
198 30 25,50 25,50 72,-30 72,-30 25))". The script produced a list of JSON files, which we
199 subsequently converted to CSV format to facilitate handling.

200

201 ***Format type and data availability***

202 The dataset file is in comma-separated values (csv) format, not compressed. Data are
203 available in Figshare at the following Digital Object Identifier:
204 <https://doi.org/10.6084/m9.figshare.30696173>.

205

206 **Header information**

207 The variables included in the dataset were selected in accordance with the Darwin Core
208 standard (Wieczorek et al., 2012), and the corresponding categories are listed in Table S1.
209 Headers are mostly self-explanatory. The dataset is fully interoperable with the European
210 subterranean spider trait dataset (<https://doi.org/10.6084/m9.figshare.16574255>), allowing
211 the extraction of morphological and ecological trait information for each species (Mammola
212 et al., 2022; Patiño-Sauma et al., 2025).

213

214 **Taxonomic validation**

215 We standardized taxonomy to the species level when feasible, following the latest
216 nomenclature of the World Spider Catalog (2025). Furthermore, we incorporated genus-
217 level records with uncertain specific attribution (e.g., *Meta* sp., *Troglohyphantes* cf.
218 *lucifuga*), as well as morphospecies under description (labelled as Genus + an
219 alphanumeric code [e.g., *Meta* sp.1]) to ensure maximal dataset breadth. Please refer to
220 the column “acceptedNameUsage” for the most up-to-date, validated taxonomic attribution
221 for each record (note that taxonomy will be updated with any new release of the dataset).
222 Eventual remarks on taxonomic decisions are provided in column “identificationRemarks”.

223

224 **Spatial validation**

225 We validated geographic coordinates (based on the WGS84 datum) through cross-
226 referencing with online resources (e.g., speleological cadastres) and, where available,
227 species-specific reference materials. We subsequently projected and visually inspected
228 localities using both R (R Core Team, 2025) and QGIS (QGIS.org, 2025). We harmonized
229 cave names and their corresponding coordinates, obtained from various sources, through
230 additional verification with national speleological cadastres whenever possible.
231 Notwithstanding these quality checks, due to missing information (especially for old
232 records), the precision of 1,990 records remains low (e.g., georeferenced using the
233 centroid of the municipality) and 108 records lack coordinates. Uncertainty in the precision
234 of coordinates is provided in the column georeferenceRemarks.

235

236 **Prediction of species richness**

237 We illustrate a potential usage of the dataset by predicting species richness patterns in
238 Europe. We achieved this by using stack species distribution modelling (SSDM) to
239 calculate species potential distribution across the continent. A ODMAP (Overview, Data,
240 Model, Assessment and Prediction) (Zurell et al., 2020) protocol for the model, detailing
241 the main analytical steps, is available in the supplementary materials (Appendix 1).

242 We included in the modelling all species with at least 10 independent records,
243 defined as occurrences from different localities separated by at least 10 km (i.e., the
244 spatial resolution of our environmental predictors), with a total of 99 different spider
245 species distributions modelled in the present study. Note that species with fewer than 10
246 independent records were only later included in the final richness prediction (see below).
247 Since some localities in the species’ distribution data have highly precise geo-localization,

multiple points may fall within the same cell of the downloaded abiotic layers, either due to the accuracy of the coordinates or because specimens were collected in nearby caves. To avoid redundant distribution points, we adjusted all coordinates to match the centroid of the corresponding cell. We then applied spatial thinning using the *thin* function from the R package “spThin” (Aiello-Lammens et al. 2015).

Previous research has shown the importance of present (Mammola and Leroy 2018) and past climatic factors (Hewitt 1999; Mammola et al. 2018, 2019b; Knüsel et al. 2024), as well as soil composition (Pavlek and Mammola 2021) in shaping subterranean species distributions. To include these variables in our model, we downloaded climatic and elevation data from the WorldClim 2 database (Fick and Hijmans 2017), specifically annual mean temperature (BIO 1), temperature seasonality (BIO 4), maximum temperature of the warmest month (BIO 5), annual precipitation (BIO 12), precipitation seasonality (BIO 15), precipitation of the warmest quarter (BIO 18) and precipitation of the coldest quarter (BIO 19). All these variables have been shown to be good proxies for subterranean climatic conditions (Mammola and Leroy 2018). In addition, we included the differences between the present and past precipitation and temperature compared during the last glacial maximum (LGM). For soil composition, we downloaded layers from the SoilGrids database (Poggio et al. 2021), specifically the percentage of coarse fragments and the gravimetric content of sand and clay in the soil. Finally, we included layers regarding evapotranspiration (Muñoz Sabater, J. (2019), normalized difference vegetation index (NDVI) from Li et al. (2023), and the soil organic carbon (SOC) and organic carbon detection (OCD) downloaded using the *soil_world* function from the “geodata” R package (Hijmans et al. 2024). All these variables are potential proxies for energy availability within the subterranean domain. All layers had a resolution of approximately 10km. We calculated pairwise Pearson’s *r* correlation coefficients among these variables, and excluded those with high correlation ($|r| \geq 0.7$), retaining only one variable from each correlated group based on ecological relevance, data quality, and consistency with previous subterranean ecology studies (Mammola & Leroy, 2018). The final list of predictors included: temperature seasonality, precipitation of the warmest quarter, evapotranspiration, percentage of coarse fragments and content of clay in the soil.

We generated individual species models using the *modelling* function from the R package “SSDM” (Schmitt et al. 2017) with the MAXENT algorithm (Phillips et al., 2004, 2006; Elith et al., 2011). To estimate species’ potential distributions and reduce overprediction, thresholds on environmental suitability were applied using Cohen's Kappa and True Skill Statistic (TSS) values via the *ecospat.max.kappa* and *ecospat.max.tss* functions from the R package “ecospat” (Broennimann et al. 2025), with the more restrictive threshold being selected. We then stacked all individual species distributions, and included species with fewer than 10 records as single-cell localities in the map.

Results

The dataset contains 31,224 records of subterranean spiders, accounting for 637 species, comprising all available georeferenced data on 31,116 records up to 2025, spanning 40 countries also including Azerbaijan, Georgia, Turkey, and Russia.

291 The spatial extent of the dataset ranges from -28.80°W to 50.02°E in longitude, and from
292 27.65°N to 67.94°N in latitude.

293 Occurrence densities are particularly high in several European regions. Northern
294 Italy, especially the northwestern and northeastern Alps, shows the greatest number of
295 records. Elevated densities also characterize northern Spain along the Atlantic coast, as
296 well as Slovenia and Croatia within the Dinaric karst (Fig. 1A). Observed species richness
297 closely matches these patterns. The highest values occur in Croatia (Dinaric karst) and
298 Slovenia, where multiple taxa overlap, followed by parts of the Alps and northern Spain's
299 Atlantic region (Fig. 1B).

300 A country-level summary highlights marked geographic disparities. Italy stands out for both
301 richness and number of records, followed by Croatia, Spain, France, and Slovenia (Fig.
302 1C). Most other European countries display comparatively low values, underscoring the
303 strong imbalance in sampling effort.

304 Predicted species richness (Fig. 1D) confirms the Alps and adjacent mountain systems,
305 including the Dinaric Arc, as major hotspots, with high values spanning northern Italy,
306 Slovenia, Croatia, and the Atlantic coast of northern Spain. Moderate richness also
307 extends into parts of central and southeastern Europe, including Austria, Germany,
308 France, and Bosnia and Herzegovina.

309

310

311

312

313

314

315

316

317

318

319

320

321

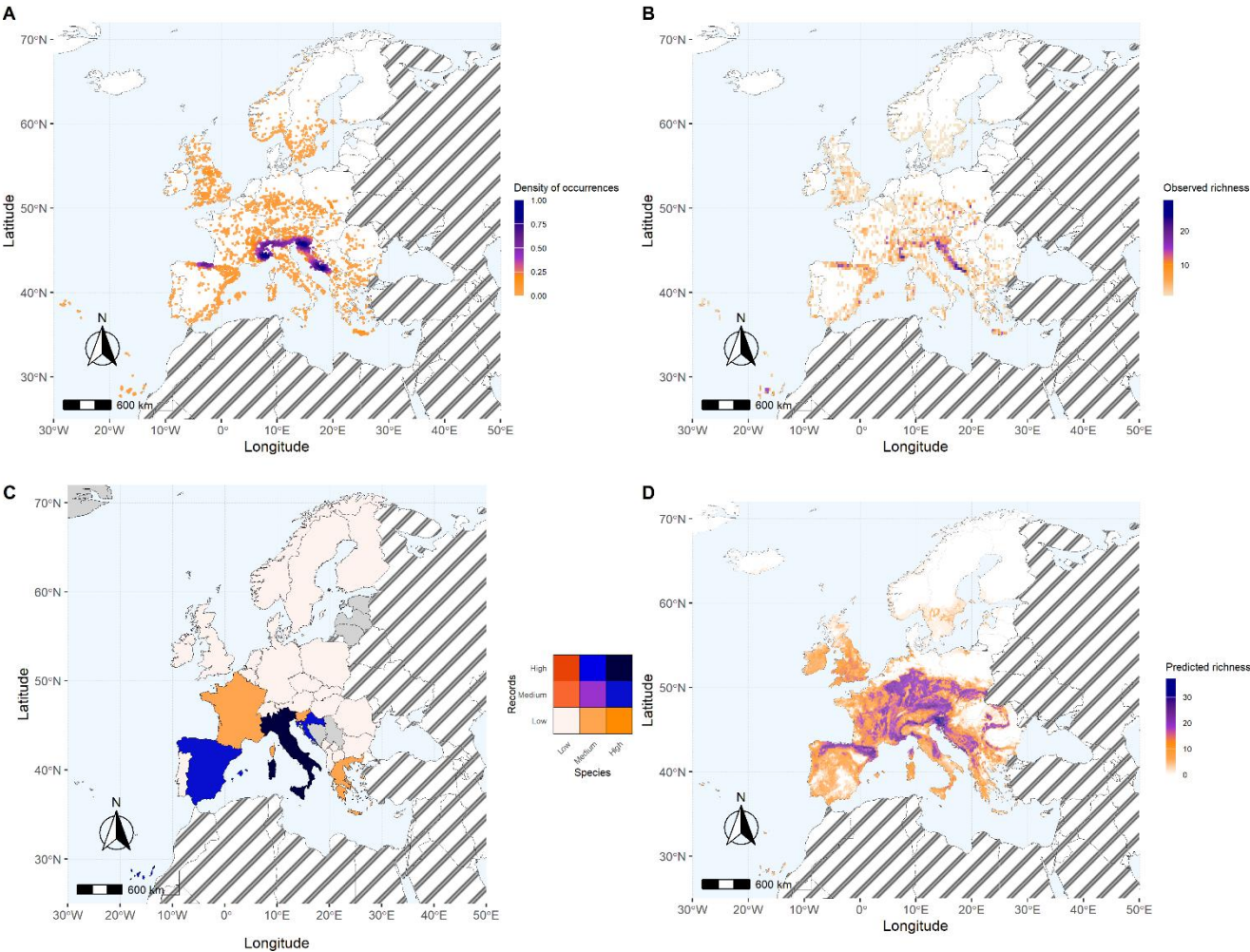
322

323

324

325

326



328

329 **Figure 1. (A)** Density of species occurrences across Europe. Points represent individual records, with color intensity reflecting local occurrence density; **(B)** Observed richness shows the number of
330 species recorded per cell across the continent. **(C)** Geographic distribution of spider records in Europe. The map shows the combined effect of the number of records and the species richness for
331 each country, classified into tertiles (Low, Medium, High). The 3×3 legend indicates the intersection between record abundance (rows) and species richness (columns). Countries without records are
332 shown in grey. **(D)** Predicted richness shows the modelled species richness across Europe obtained through stacked species distribution modelling, with colours indicating the predicted
333 number of species per grid cell. Mapping is restricted to countries with sufficient data availability; countries with very few records (e.g. Russia) are therefore not represented.
334
335
336
337
338

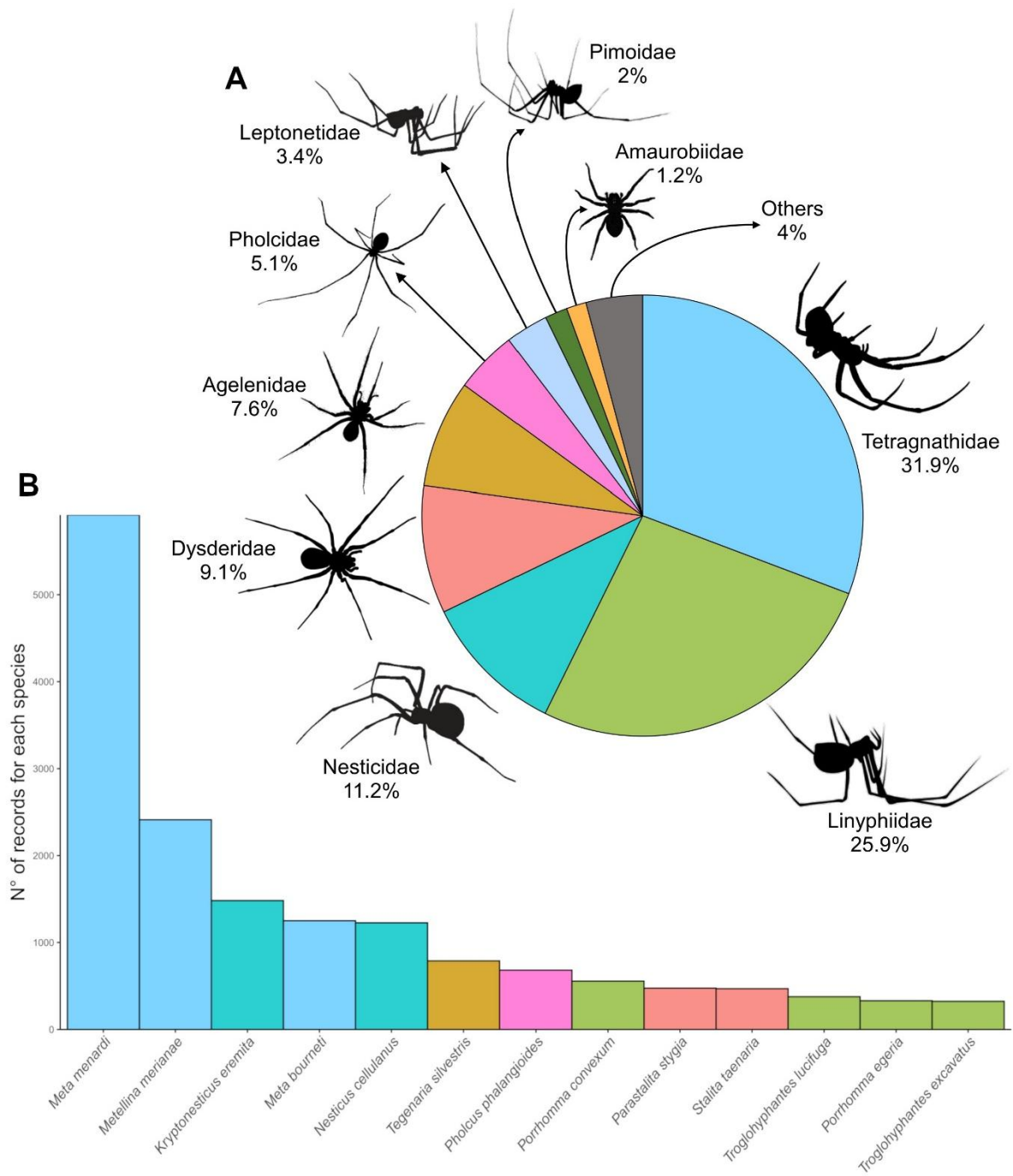
339

340 Most records were collected from natural caves, totalling 25,947 (83.10%) of all
341 occurrences. These were followed by artificial habitats (e.g., mines, bunkers, blockhouses,
342 cellars), which accounted for 1,929 records (6.17%), and shallow subterranean habitats
343 (SSH) with 323 records (1.03%).

344 In terms of taxonomic distribution, most records belong to the family Tetragnathidae
345 (9,951; 31.87%), followed by Linyphiidae (8,094; 25.93%), Nesticidae (3,483; 11.16%),
346 Dysderidae (2,822; 9.04%), Agelenidae (2,378; 7.62%), Pholcidae (1,584; 5.07%),
347 Leptonetidae (1,057; 3.39%), Pimoidae (618; 1.98%), and Amaurobiidae (362; 1.16%).

348 Other families, including Sicariidae, Theridiidae, and Cicurinidae, were each represented
349 by less than 1% of the total records (Fig. 2A).

350



351

352 **Figure 2. (A)** Pie chart showing the relative percentage of occurrences records for each spider
353 family in the dataset. Families representing less than 1% of records were grouped into the category
354 “Others”. **(B)** Barplot displaying the number of records for the most common species in the spiders
355 dataset. Only species with at least 300 records are shown.

356

357 Among the 637 species in the dataset, the most frequently recorded species is *Meta*
358 *menardi* (Latreille, 1804), with 5,914 occurrences. This is followed by *Metellina merianae*
359 (Scopoli, 1763) with 2,412 records, *Kryptonesticus eremita* (Simon, 1880) with 1,482
360 records, *Meta bourneti* (Simon, 1922) with 1,251 records, and *Nesticus cellulanus* (Clerck,
361 1757) with 1,226 records. All other species are represented by fewer than 1,000 records
362 (average number of records \pm s.d. = 20 ± 38) (Fig. 2B).

363

364 Discussion

365 The dataset provides the most comprehensive distribution data for subterranean spiders
366 across Europe. It brings together data from multiple countries and research groups,
367 consolidating previously scattered information into a single, standardized, and harmonized
368 resource.

369 The taxonomic composition of the dataset reveals meaningful ecological and
370 sampling-related patterns. In particular, some families are overrepresented in terms of
371 occurrence records relative to their species richness. For example, Tetragnathidae account
372 for a high number of records despite including few subterranean-associated species (*Meta*
373 spp. and *Metellina merianae*; Fig. 2A). This likely reflects the large body size, high
374 conspicuousness, and often high local densities of these spiders near cave entrances
375 (Smithers 2005; Novak et al., 2010; Mammola & Isaia, 2014), all resulting in high
376 detectability, as well as their broad geographic distributions in Europe (Mammola et al.,
377 2019a, 2021c).

378 The predicted richness patterns (Fig. 1D) also illustrate the potential of the dataset to
379 support macroecological inference. By integrating species distribution models for tens of
380 taxa, the resulting map highlights broad-scale biogeographic structures that are not always
381 evident from raw data alone. In particular, the Alps, the Dinaric Arc, and parts of northern
382 Spain emerge as major hotspots of subterranean diversity, consistent with the long-term
383 persistence of stable microclimatic refugia and the complex geomorphological history of
384 these regions (Culver & Sket, 2000; Deharveng et al., 2024). The smoother gradients
385 revealed by the prediction also indicate that true richness likely extends beyond areas with
386 dense sampling, highlighting regions such as Austria, Germany, France, and Bosnia and
387 Herzegovina as potentially important, yet comparatively understudied. This underscores
388 both the ecological value of the modelling approach and the role of the dataset in identifying
389 priority areas where additional sampling would substantially improve knowledge of
390 subterranean biodiversity. Integrating predicted patterns with conservation planning may
391 therefore help guide efforts toward mountain systems and karstic landscapes that harbour
392 disproportionate levels of subterranean diversity.

393 The initiative reflects the growing recognition that open data sharing at continental
394 and global scales is essential to advance (macro)ecological research, improve biodiversity
395 monitoring, and inform conservation strategies under increasing environmental threats
396 (Urbano et al., 2024). This goal has been greatly facilitated in recent years by the
397 development of international research infrastructures such as the Global Biodiversity
398 Information Facility (GBIF; GBIF.org, 2025) and LifeWatch ERIC (Basset & Los, 2012), which
399 promote standardized, interoperable, and openly accessible biodiversity data.

Although focused on distributional records, its interoperability with other data sources makes it particularly valuable. For example, it can be combined with complementary resources, such as datasets on spider morphological and ecological traits (Mammola et al., 2022; Patiño-Sauma et al., 2025) or phylogenetic information from other sources, thereby enabling comparative analyses, functional diversity assessments, and large-scale ecological modelling. By linking distributional data with traits and phylogenies, researchers can further explore questions related to the processes shaping subterranean biodiversity and identify species and regions most vulnerable to environmental change and anthropogenic pressures.

The geographical and taxonomic breadth of the dataset makes it a valuable resource for addressing key questions in subterranean ecology, from species distributions and environmental drivers to long-standing gaps in conservation status, functional diversity, and the ecological factors structuring subterranean spider assemblages. While the dataset marks a substantial step forward, some regions and taxa remain underrepresented, reflecting historical biases in research effort and data availability. Open data sharing, however, provides the basis for progressively improving coverage and quality, fostering new sampling initiatives, comparative analyses, and multi-scale investigations of subterranean biodiversity.

References

- Ahyong, S., Boyko, C. B., Bernot, J., et al. (2025). World Register of Marine Species. VLIZ. Available from: marinespecies.org (accessed 23 August 2025).
<https://doi.org/10.14284/170>
- Aiello-Lammens, M. E., Boria, R. A., Radosavljevic, A., Vilela, B. & Anderson, R. P. (2015). spThin: an R package for spatial thinning of species occurrence records for use in ecological niche models. *Ecography*, 38, 541–545.
- Bai, X., Zhang, P., Gan, L., et al. (2025). Global diversity patterns and threats of cave fish. *Global Ecology and Biogeography*, 34, e70160.
- Basset, A. & Los, W. (2012). Biodiversity e-science: LifeWatch, the European infrastructure on biodiversity and ecosystem research. *Plant Biosystems*, 146, 780–782.
<https://doi.org/10.1080/11263504.2012.740091>
- Broennimann, O., Di Cola, V. & Guisan, A. (2025). ecospat: Spatial ecology miscellaneous methods. R package version 4.1–2.
- Cancellario, T., Golomb Duran, T., Far Morenilla, A. J., Roldan, A. & Capa, M. (2025). biodumpy: A comprehensive biological data downloader. *bioRxiv*.
- Cardoso, P., Erwin, T. L., Borges, P. A. V. & New, T. R. (2011). The seven impediments in invertebrate conservation and how to overcome them. *Biological Conservation*, 144, 2647–2655.
- Cardoso, P., Pekár, S., Birkhofer, K., Chuang, A., Fukushima, C. S., Hebets, E. A., ... Mammola, S. (2025). Ecosystem services provided by spiders. *Biological Reviews*.

440 Cerasoli, F., Fiasca, B., Di Cicco, et al. (2025). EGCop: An expert-curated occurrence
441 dataset of European groundwater-dwelling copepods. *Global Ecology and*
442 *Biogeography*, 34, e13953.

443 Crespo, L. C., Pereira, F., Amorim, I. R. & Borges, P. A. V. (2025). Insights from the
444 Dalberto Teixeira Pombo (DTP) Arthropod Collection – I. Revealing the hidden
445 diversity of terrestrial cave arthropods in the Azores. *Biodiversity Data Journal*, 13,
446 e158467. <https://doi.org/10.3897/BDJ.13.e158467>

447 Culver, D. C. & Sket, B. (2000). Hotspots of subterranean biodiversity in caves and wells.
448 *Journal of Cave and Karst Studies*, 62, 11–17.

449 De Biurrun, G., Prieto, C. & Baquero, E. (2022). ArachnoMap, una herramienta para
450 difundir el conocimiento del taxón Araneae en la Península Ibérica y Baleares.
451 *Revista Ibérica de Aracnología*, 40, 2–3.

452 Deharveng, L., Stoch, F., Gibert, J., Bedos, A., Galassi, D., Zagmajster, M., ... Marmonier,
453 P. (2009). Groundwater biodiversity in Europe. *Freshwater Biology*, 54, 709–726.
454 <https://doi.org/10.1111/j.1365-2427.2008.01972.x>

455 Deharveng, L., Bedos, A., Pipan, T. & Culver, D. C. (2024). Global subterranean
456 biodiversity: A unique pattern. *Diversity*, 16, 157.

457 Dornelas, M., Antão, L. H., Bates, A. E., Brambilla, V., Chase, J. M., Chow, C. F., ... Fryxell,
458 J. (2025). BioTIME 2.0: Expanding and improving a database of biodiversity time
459 series. *Global Ecology and Biogeography*, 34, e70003.

460 Elith, J., Phillips, S. J., Hastie, T., Dudík, M., Chee, Y. E. & Yates, C. J. (2011). A statistical
461 explanation of MaxEnt for ecologists. *Diversity and Distributions*, 17, 43–57.
462 <https://doi.org/10.1111/j.1472-4642.2010.00725.x>

463 Fialas, P. C., Santini, L., Russo, D., et al. (2025). Changes in community composition and
464 functional diversity of European bats under climate change. *Conservation Biology*,
465 39, e70025.

466 Ficetola, G. F., Canedoli, C. & Stoch, F. (2019). The Racovitza impediment and the
467 hidden biodiversity of unexplored environments. *Conservation Biology*, 33, 214–
468 216.

469 Fick, S. E. & Hijmans, R. J. (2017). WorldClim 2: new 1 km spatial resolution climate
470 surfaces for global land areas. *International Journal of Climatology*, 37, 4302–4315.

471 GBIF: The Global Biodiversity Information Facility (2025). What is GBIF? Available from:
472 <https://www.gbif.org/what-is-gbif>

473 Gobierno de Canarias (2025). Banco de Datos de Biodiversidad de Canarias. Available
474 from: <https://www.biodiversidadcanarias.es/biota>
475 (accessed 30 May 2025).

476 Hewitt, G. M. (1999). Post-glacial re-colonization of European biota. *Biological Journal of*
477 *the Linnean Society*, 68, 87–112.

478 Hijmans, R. J., Barbosa, M., Ghosh, A. & Mandel, A. (2024). geodata: Download
 479 geographic data. R package version 0.6-2. Available from: [https://CRAN.R-](https://CRAN.R-project.org/package=geodata)
 480 [project.org/package=geodata](https://CRAN.R-project.org/package=geodata)

481 Hortal, J., de Bello, F., Diniz-Filho, J. A. F., Lewinsohn, T. M., Lobo, J. M. & Ladle, R. J.
 482 (2015). Seven shortfalls that beset large-scale knowledge of biodiversity. *Annual*
 483 *Review of Ecology, Evolution, and Systematics*, 46, 523–549.

484 Hughes, A. C., Orr, M. C., Ma, K., Costello, M. J., Waller, J., Provoost, P., ... Qiao, H.
 485 (2021). Sampling biases shape our view of the natural world. *Ecography*, 44, 1259–
 486 1269.

487 Keith, D. A., Ferrer-Paris, J. R., Nicholson, E. & Kingsford, R. T. (eds.) (2020). The IUCN
 488 global ecosystem typology 2.0. Gland, IUCN.

489 Keith, D. A., Ferrer-Paris, J. R., Nicholson, E., et al. (2022). A function-based typology for
 490 Earth's ecosystems. *Nature*, 610, 513–518.

491 Knüsel, M., Alther, R. & Altermatt, F. (2024). Pronounced changes of subterranean
 492 biodiversity patterns along a Late Pleistocene glaciation gradient. *Ecography*,
 493 e07321.

494 Labouyrie, M., Ballabio, C., Romero, et al. (2023). Patterns in soil microbial diversity
 495 across Europe. *Nature Communications*, 14, 3311.

496 LifeWatch ERIC (2020). EcoPortal, the repository of semantic resources for the ecological
 497 domain. Available from: <https://ecoportal.lifewatch.eu>

498 Macías-Hernández, N., Suárez, D., de la Cruz-López, S., López, H. & Oromí, P. (2024).
 499 Diversidad de arañas hipogeas del archipiélago canario. *Ecosistemas*, 33, 2516.

500 Mammola, S., Giachino, P. M., Piano, E., Jones, A., Barberis, M., Badino, G. & Isaia, M.
 501 (2016). Ecology and sampling techniques of the Milieu Souterrain Superficiel
 502 (MSS). *The Science of Nature*, 103, 88.

503 Mammola, S. & Isaia, M. (2014). Niche differentiation in *Meta bourneti* and *M. menardi*
 504 (Araneae, Tetragnathidae) with notes on the life history. *International Journal of*
 505 *Speleology*, 43, 343–353.

506 Mammola, S. & Isaia, M. (2017). Spiders in caves. *Proceedings of the Royal Society B:*
 507 *Biological Sciences*, 284, 20170193.

508 Mammola, S., Goodacre, S. L. & Isaia, M. (2018). Climate change may drive cave spiders
 509 to extinction. *Ecography*, 41, 233–243.

510 Mammola, S. & Leroy, B. (2018). Applying species distribution models to caves and other
 511 subterranean habitats. *Ecography*, 41, 1194–1208.

512 Mammola, S., Schönhöfer, A. L. & Isaia, M. (2019). Tracking the ice: Subterranean
 513 harvestmen distribution matches ancient glacier margins. *Journal of Zoological*
 514 *Systematics and Evolutionary Research*, 57, 548–554.

515 Mammola, S., Cardoso, P., Angyal, D., et al. (2019). Continental data on cave-dwelling
 516 spider communities across Europe. *Biodiversity Data Journal*, 7, e38492.

517 Mammola, S., Souza, M. F. V. R., Isaia, M. & Ferreira, R. L. (2021). Global distribution of
518 microwhip scorpions. *Journal of Biogeography*, 48, 1518–1529.

519 Mammola, S., Lunghi, E., Bilandžija, H., Cardoso, P., Grimm, V., Schmidt, S. I., ...
520 Martínez, A. (2021). Collecting eco-evolutionary data in the dark. *Ecology and*
521 *Evolution*, 11, 5911–5926.

522 Mammola, S., Hesselberg, T. & Lunghi, E. (2021). A trade-off between latitude and
523 elevation contributes to explain range segregation of broadly distributed cave-
524 dwelling spiders. *Journal of Zoological Systematics and Evolutionary Research*, 59,
525 370–375.

526 Mammola, S., Pavlek, M., Huber, B. A., et al. (2022). A trait database and updated
527 checklist for European subterranean spiders. *Scientific Data*, 9, 236.

528 Martínez, A., Dragomir-Cosmin, D., Cancellario, T., et al. (2024). Stygofauna Mundi: A
529 comprehensive global biodiversity database of groundwater related habitats across
530 marine and freshwater realms. Proceedings of the 26th International Conference on
531 Subterranean Biology and 6th International Symposium on Anchialine Ecosystems.
532 SUPSI.

533 Martínez-Núñez, C., Martínez-Prentice, R. & García-Navas, V. (2023). Land-use diversity
534 predicts regional bird taxonomic and functional richness worldwide. *Nature*
535 *Communications*, 14, 1320.

536 Mori, N., Vehovar, Ž., Brad, T., et al. (2025). A comprehensive occurrence dataset for
537 European Ostracoda. *Global Ecology and Biogeography*, 34, e70065.

538 Muñoz Sabater, J. (2019). ERA5-Land monthly averaged data from 1950 to present.
539 Copernicus Climate Change Service (C3S) Climate Data Store (CDS).
540 <https://doi.org/10.24381/cds.68d2bb30>
541 (accessed 30 June 2025).

542 Li, M., Cao, S., Zhu, Z., Wang, Z., Myneni, R. B. & Piao, S. (2023). Spatiotemporally
543 consistent global dataset of the GIMMS normalized difference vegetation index
544 (PKU GIMMS NDVI) from 1982 to 2022 (V1.2) [Data set]. Zenodo.
545 <https://doi.org/10.5281/zenodo.8253971>

546 Nicolosi, G., Martínez García, A., Piano, E., Isaia, M. & Mammola, S. (2025). An expert-
547 curated dataset on cave-dwelling spider communities in the Western Italian Alps.
548 *Biogeographia*, 40, 1–17.

549 Novak, T., Tkvac, T., Kuntner, M., Arnett, E. A., Delakorda, S. L., Perc, M. & Janžekovič, F.
550 (2010). Niche partitioning in orbweaving spider *Meta menardi* and *Metellina*
551 *merianae* (Tetragnathidae). *Acta Oecologica*, 36, 522–529.

552 Pantini, P. & Isaia, M. (2019). Araneae.it: the online catalog of Italian spiders. *Fragmenta*
553 *Entomologica*, 51, 127–152.

554 Paragamian, K., Poulinakis, M., Paragkamian, S. & Nikoloudakis, I. (2025). Cave fauna of
555 Greece database – Hellenic Institute of Speleological Research. Available from:
556 <https://database.inspee.gr> (accessed 26 July 2025).

557 Patiño-Sauma, D., Cardoso, P., Oromí, P., et al. (2025). Another brick in the wall of
558 European subterranean spider knowledge: adding Macaronesian species. *bioRxiv*.

559 Pavlek, M. & Mammola, S. (2021). Niche-based processes explaining the distributions of
560 closely related subterranean spiders. *Journal of Biogeography*, 48, 118–133.

561 Poggio, L., De Sousa, L. M., Batjes, N. H., Heuvelink, G. B. M., Kempen, B., Ribeiro, E. &
562 Rossiter, D. (2021). SoilGrids 2.0. *SOIL*, 7, 217–240.

563 Phillips, S. J., Anderson, R. P. & Schapire, R. E. (2006). Maximum entropy modeling of
564 species geographic distributions. *Ecological Modelling*, 190, 231–259.

565 Phillips, S. J., Dudík, M. & Schapire, R. E. (2004). A maximum entropy approach to
566 species distribution modeling. In *Proceedings of the Twenty-First International*
567 *Conference on Machine Learning (ICML)*, 472–486.

568 QGIS.org (2025). QGIS geographic information system. QGIS Association.

569 R Core Team (2025). R: A language and environment for statistical computing. Vienna,
570 Austria, R Foundation for Statistical Computing.

571 Sabatini, F. M., Jiménez-Alfaro, B., Jandt, U., Chytrý, M., Field, R., Kessler, M., ...
572 Bruehlheide, H. (2022). Global patterns of vascular plant alpha diversity. *Nature*
573 *Communications*, 13, 4683.

574 Saclier, N., Duchemin, L., Konecny-Dupré, L., et al. (2024). The World Asellidae database.
575 *Molecular Ecology Resources*, 24, e13882.

576 Schmitt, S., Pouteau, R., Justeau, D., De Boissieu, F. & Birnbaum, P. (2017). ssdm: An R
577 package to predict species richness and composition. *Methods in Ecology and*
578 *Evolution*, 8, 1795–1803.

579 Smithers, P. (2005). The early life history and dispersal of the cave spider *Meta menardi*
580 (Latreille, 1804) (Araneae: Tetragnathidae). *Bulletin of the British Arachnological*
581 *Society*, 13, 213–216.

582 Tanalgo, K. C., Tabora, J. A. G., de Oliveira, H. F. M., et al. (2022). DarkCideS 1.0: A global
583 database for bats in karsts and caves. *Scientific Data*, 9, 155.

584 Turnbull, A. L. (1973). Ecology of the true spiders (Araneomorphae). *Annual Review of*
585 *Entomology*, 18, 305–348.

586 Urbano, F., Viterbi, R., Pedrotti, L., Vettorazzo, E., Movalli, C. & Corlatti, L. (2024).
587 Enhancing biodiversity conservation and monitoring in protected areas through
588 efficient data management. *Environmental Monitoring and Assessment*, 196, 12.

589 Wieczorek, J., Bloom, D., Guralnick, R., Blum, S., Döring, M., et al. (2012). Darwin Core:
590 An evolving community-developed biodiversity data standard. *PLOS ONE*, 7,
591 e29715. <https://doi.org/10.1371/journal.pone.0029715>

592 World Spider Catalog (2025). World Spider Catalog, version 26. Natural History Museum
593 Bern. Available from: <http://wsc.nmbe.ch> (accessed 31 August 2025).
594 <https://doi.org/10.24436/2>.

- 595 Zagmajster, M., Eme, D., Fišer, C., Galassi, D., Marmonier, P., Stoch, F. & Malard, F.
 596 (2014). Geographic variation in range size and beta diversity of groundwater
 597 crustaceans: Insights from habitats with low thermal seasonality. *Global Ecology*
 598 *and Biogeography*, 23, 1135–1145. <https://doi.org/10.1111/geb.12200>
- 599 Zagmajster, M. (2016). SubBioDatabase – a tool for research and conservation of
 600 subterranean biodiversity of the whole Dinarides. In M. Lukić (Ed.), Abstract book.
 601 Zagreb, Croatian Biospeleological Society, p. 41.
- 602 Zurell, D., Franklin, J., König, C., Bouchet, P. J., Serra-Diaz, J. M., Dormann, C. F., ...
 603 Merow, C. (2020). A standard protocol for describing species distribution models.
 604 *Ecography*, 43, 1261–1277. <https://doi.org/10.1111/ecog.04960>

605

606 **Funding**

607 This research was funded by Biodiversa+ (project ‘DarCo’), the European Biodiversity
 608 Partnership under the 2021–2022 BiodivProtect joint call for research proposals, co-
 609 funded by the European Commission (GA N°101052342) and with the funding
 610 organizations Ministry of Universities and Research (Italy), Agencia Estatal de
 611 Investigación—Fundación Biodiversidad (Spain), Fundo Regional para a Ciência e
 612 Tecnologia (Portugal), Suomen Akatemia—Ministry of the Environment (Finland), Belgian
 613 Science Policy Office (Belgium), Agence Nationale de la Recherche (France), Deutsche
 614 Forschungsgemeinschaft e.V. (Germany), Schweizerischer Nationalfonds (Grant No.
 615 31BD30_209583, Switzerland), Fonds zur Förderung der Wissenschaftlichen Forschung
 616 (Austria), Ministry of Higher Education, Science and Innovation (Slovenia), Ministry of
 617 Agriculture of the Czech Republic (MZe RO0425, Czechia) and the Executive Agency for
 618 Higher Education, Research, Development and Innovation Funding (Romania) and
 619 Biodiversa+, the European Biodiversity Partnership, in the context of the Sub-BioMon -
 620 Developing and testing approaches to monitor subterranean biodiversity in karst project
 621 under the 2022-2023 BiodivMon joint call. It was co-funded by the European Commission
 622 (GA N°101052342) and the following funding organisations: Ministry of Higher Education,
 623 Science and Innovation (Slovenia), The Belgian Science Policy (Belgium), Ministry of
 624 Universities and Research (Italy), National Research, Development and Innovation Office
 625 (Hungary), Executive Agency for Higher Education, Research, Development and
 626 Innovation Funding (Romania), and self-financing partner National Museum of Natural
 627 History Luxembourg (Luxembourg). Additional support to S.M. is provided by the P.R.I.N.
 628 2022 “DEEP CHANGE” (2022MJSYF8), funded by the Ministry of Universities and
 629 Research (Italy), and by NBFC, funded by the Italian Ministry of University and Research,
 630 P.N.R.R., Missione 4, Componente 2, “Dalla ricerca all’impresa”, Investimento 1.4, Project
 631 CN00000033. Support to T.D. and M.Z. is provided by Slovenian Agency for Research and
 632 Innovation via programme support (programme P1-0184).

633

634 **Conflict of Interest Statement**

635 The authors declare no conflicts of interest.

636

637 **Data and Code Availability Statement**

638 The dataset and associated R codes are available from the Figshare repository under
639 <https://doi.org/10.6084/m9.figshare.30696173>

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668 **Appendix 1.** Summary of the modelling pipeline according to the
669 ODMAP (Zurell et al., 2020) reporting protocol.

672 **OVERVIEW**

673 **Title**

674 An open occurrence dataset for European subterranean spider

675 **Model objective**

676 **Model objective:** Mapping and interpolation.

677 **Target output:** Obtaining a map of species richness based on stacked predicted distribution ranges of each
678 species.

679 **Focal Taxon**

680 **Focal Taxon:** Spiders (Arachnida: Araneae).

681 **Location**

682 **Location:** Europe.

683 **Scale of Analysis**

Spatial extent: -30, 50, 25, 72 (xmin, xmax, ymin, ymax).
Spatial resolution: 10 km.
Temporal extent: Present.
Temporal resolution: NA.
Boundary: political.

685 **Biodiversity data**

Observation type: field survey, GPS tracking.
Response data type: presence-only.

687 **Predictors**

688 **Predictor types:** climatic, habitat, edaphic.

689 **Hypotheses**

690 **Hypotheses:** No specific hypotheses were made regarding the species-environment relationships.

691 **Assumptions**

692 **Model assumptions:** Unlimited dispersal.

693 **Algorithms**

Modelling techniques:	maxent.
Model complexity:	Only ecologically interpretable predictors were included, and interactions were not considered to maintain biological interpretability.
Model averaging:	NA.

694

695 **Workflow**

696 Environmental predictor layers were standardized prior to modelling.

697 For each species, a MaxEnt model was fitted using presence locations and environmental predictors to
698 generate continuous habitat suitability predictions.

699 Predictions were spatially constrained using a buffer-based mask derived from the mean and maximum
700 inter-point distances among occurrence records, and suitability values were smoothed as a function of
701 distance to the accessible area boundary.

702 Model performance was assessed by extracting predicted suitability values at observed presences and
703 randomly sampled background points, and optimal thresholds were determined using both Kappa and True
704 Skill Statistic (TSS).

705 Binary presence–absence maps were generated using the most conservative threshold and combined
706 across species to produce a stacked species distribution model representing predicted species richness.

707 **Software**

708 **Software:** R version 4.5.0. SSDM package version 0.2.9.

Code availability: The code to run this analysis is available in the Figshare repository (<https://doi.org/10.6084/m9.figshare.30696173>).
Data availability: Distribution data are available in the Figshare repository (<https://doi.org/10.6084/m9.figshare.30696173>). Environmental predictors are available from different sources (details in section: Predictor variables).

709 **DATA**

710 **Biodiversity data**

711 **Taxon names:** 99 different spider species (including some under description) were modelled. Taxonomy is
712 according to the World Spider Catalog version 26 (<https://wsc.nmbe.ch/>) and the version 3 of the Checklist
713 of European Subterranean spider (<https://doi.org/10.6084/m9.figshare.16574255>). See supplementary
714 material of the manuscript for a detailed list of all species analyzed.

715 **Taxonomic reference system**

716 **Ecological level:** species.

717 **Data sources**

718 **Sampling design:** Random.

719 **Sample size:** 7000 different observations for 99 species (average number of records per species \pm s.d.: 70.7
720 \pm 180.7).

721 **Clipping:** Europe.

722 **Scaling:** We used one record per cell for each species.

723 **Cleaning:** Species with less than 10 records were discarded for the modelling analyses. For species with
724 less than 10 records, the observed distribution was used to generate the stacked prediction of species
725 richness. To address spatial sampling bias and spatial autocorrelation, occurrence records were spatially
726 thinned separately for each species using a minimum nearest-neighbour distance of 10 km. The thinning
727 procedure was repeated 100 times and a single thinned realization was retained for model fitting.
728 Environmental predictor layers and model predictions were processed on a common raster grid and spatial
729 resolution to ensure spatial consistency in subsequent thresholding and stacking.

730 **Absence data:** No true absence data were available.

731 **Background data:** Background data were derived separately for each species by defining a species-specific
732 accessible area based on the spatial configuration of occurrence records. The mean and maximum pairwise
733 distances among occurrence locations were used to generate a buffered spatial mask representing the area
734 available for dispersal. Random background points were then sampled within this mask at approximately
735 twice the number of presence records and used for threshold selection. This spatially constrained
736 background sampling reduced the influence of environmentally unrealistic or geographically inaccessible
737 areas on model assessment.

738 **Errors and biases:** See main text for discussion.

739 **Data partitioning**

Training data:	NA.
Validation data:	NA.
Test data:	Expert-based assessment.

740

741 **Predictor variables**

742 **Predictor variables:** Temperature seasonality, precipitation of the warmest quarter, evapotranspiration,
743 percentage of coarse fragments and content of clay in the soil.

744 **Data sources:** Evapotranspiration: [https://cds.climate.copernicus.eu/datasets/reanalysis-era5-land-](https://cds.climate.copernicus.eu/datasets/reanalysis-era5-land-monthly-means?tab=overview)
745 [monthly-means?tab=overview](https://cds.climate.copernicus.eu/datasets/reanalysis-era5-land-monthly-means?tab=overview); Bioclimatic data: <https://www.worldclim.org/data/worldclim21.html>; Soil
746 data: <https://soilgrids.org/>.

747 **Spatial extent:** -30, 50, 25, 72 (xmin, xmax, ymin, ymax).

748 **Spatial resolution:** res: 0.08333333 (around 10km at the equator).

749 **Coordinate reference system:** WGS84 decimal degrees.

750 **Temporal extent:** 1970–2000 for Bioclimatic variables, 2020 for soil data, and 1950–2025 for

751 evapotranspiration.

752 **Temporal resolution:** NA.

753 **Data processing:** Environmental predictor layers were standardized using z-score normalization (mean-

754 centering and scaling by standard deviation) prior to model fitting to ensure comparability among variables.

755 **Errors and biases:** NA.

756 **Dimension reduction:** Expert-based assessment.

757 **Transfer data**

Data sources:	NA.
Spatial extent:	
Spatial resolution:	NA.
Temporal extent:	NA.
Temporal resolution:	NA.
Models and scenarios:	NA.
Data processing:	NA.
Quantification of Novelty:	NA.

758 **MODEL**

759 **Variable pre-selection**

760 **Variable pre-selection:** NA.

761 **Multicollinearity**

762 **Multicollinearity:** We tested the correlation between all variables by calculating pairwise Pearson’s r

763 correlations, setting a threshold for collinearity at $|r| > 0.7$, and then selected the variables with higher

764 ecological meaning based on expert based assessment.

765 **Model settings**

766 **maxent:** featureSet (Default), featureRule (Default), regularizationMultiplierSet (1), regularizationRule

767 (Default), convergenceThresholdSet (Default (0.00001)), samplingBiasRule (NA), samplingBiasNotes (NA),

768 targetGroupSampleSize (NA), offsetSet (NA), offsetRule (NA), expertMapProbSet (NA), expertMapProbRule

769 (NA), expertMapRateSet (NA), expertMapRateRule (NA), expertMapSkewSet (NA), expertMapSkewRule

770 (NA), expertMapShiftSet (NA), expertMapShiftRule (NA), notes (NA).

771 **Model settings (extrapolation):** Default.

772 **Model estimates**

773 **Coefficients:** Default.

774 **Parameter uncertainty:** Parameter uncertainty was not explicitly quantified in this workflow.

775 **Variable importance:** NA.

776 **Model selection - model averaging - ensembles**

777 **Model selection:** No applicable due to single model used.

778 **Model averaging:** NA.

779 **Model ensembles:** NA.

780 **Analysis and Correction of non-independence**

781 **Spatial autocorrelation:** NA.

782 **Temporal autocorrelation:** NA.

783 **Nested data:** NA.

784 **Threshold selection**

785 **Threshold selection:** Continuous habitat suitability predictions from MaxEnt were converted to binary
786 presence–absence maps using species-specific thresholds. Thresholds were determined by maximizing both
787 Cohen’s Kappa and the True Skill Statistic (TSS) based on model predictions at observed presences and
788 randomly sampled background points. The most conservative threshold among the two was applied to
789 generate final binary maps.

790 **ASSESSMENT**

791 **Performance statistics**

792 **Performance on training data:** Expert-based assessment on the stacked richness map as we were not
793 interested in single species distributions performance.

794 **Performance on validation data:** NA.

795 **Performance on test data:** NA.

796 **Plausibility check**

797 **Response shapes:** NA..

798 **Expert judgement:** Yes

799 **PREDICTION**

800 **Prediction output**

801 **Prediction unit:** Same as extent.

802 **Post-processing:** NA.

803 **Uncertainty quantification**

804 **Algorithmic uncertainty:** NA.

805 **Input data uncertainty:** NA.

806 **Parameter uncertainty:** NA.

807 **Scenario uncertainty:** NA.

808 **Novel environments:** NA.