



HAL
open science

A framework to bridge scales in distribution modeling of soil microbiota

Jonas Lembrechts, Luke Broeders, Johan De gruyter, Dajana Radujković, Irene Ramirez-Rojas, Jonathan Roger Michel Henri Lenoir, Erik Verbruggen

► To cite this version:

Jonas Lembrechts, Luke Broeders, Johan De gruyter, Dajana Radujković, Irene Ramirez-Rojas, et al.. A framework to bridge scales in distribution modeling of soil microbiota. *FEMS Microbiology Ecology*, 2020, 96 (5), <10.1093/femsec/fiaa051>. <hal-03003203>

HAL Id: hal-03003203

<https://hal.science/hal-03003203v1>

Submitted on 13 Nov 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

1 **A framework to bridge scales in distribution modelling of soil microbiota**

2

3 *Lembrechts JJ^{(1),1}, Broeders L⁽¹⁾, De Gruyter J⁽¹⁾, Radujković D⁽¹⁾, Ramirez-Rojas I⁽¹⁾, Lenoir J⁽²⁾, Verbruggen*
4 *E⁽¹⁾*

5

6 *(1) Plant and Ecosystems Research Group (PLECO), University of Antwerp, Belgium.*

7 *(2) UR 'Ecologie et Dynamique des Systèmes Anthropisées' (EDYSAN, UMR 7058*
8 *CNRS-UPJV), Université de Picardie Jules Verne, Amiens, France.*

9

10 **Corresponding author. jonas.lembrechts@uantwerpen.be*

11

12 *Orcid ID JLL: orcid.org/0000-0002-1933-0750.*

13 *Orcid ID JL: orcid.org/0000-0003-0638-9582.*

14 *Orcid ID EV: orcid.org/0000-0001-7015-1515.*

15

16

17

18 **Keywords:** *species distribution models, soil microbes, biotic interactions, joint distribution models, niche*
19 *modelling, microclimate*

20

21

¹ Corresponding author. Jonas Lembrechts, CDE C.007 Universiteitsplein 1, University of Antwerp, 2610 Wilrijk, Belgium, +3232651727, jonas.lembrechts@uantwerpen.be.

22 **Abstract**

23

24 Creating accurate habitat suitability and distribution models (HSDMs) for soil microbiota is far
25 more challenging than for aboveground organism groups. In this perspective paper, we propose
26 a conceptual framework that addresses several of the critical issues holding back further
27 applications. Most importantly, we tackle the mismatch between the broad-scale, long-term
28 averages of environmental variables traditionally used, and the environment as experienced by
29 soil microbiota themselves. We suggest using nested sampling designs across environmental
30 gradients and objectively integrating spatially hierarchic heterogeneity as covariates in HSDMs.
31 Secondly, to incorporate the crucial role of taxa co-occurrence as driver of soil microbial
32 distributions, we promote the use of joint species distribution models, a class of models that jointly
33 analyze multiple species' distributions, quantifying both species-specific environmental responses
34 (i.e. the environmental niche) and covariance among species (i.e. biotic interactions). Our
35 approach allows incorporating the environmental niche and its associated distribution across
36 multiple spatial scales. The proposed framework facilitates the inclusion of the true relationships
37 between soil organisms and their abiotic and biotic environment in distribution models, which is
38 crucial to improve predictions of soil microbial redistributions as a result of global change.

39

40 **Introduction**

41

42 Habitat suitability and distribution models (HSDMs) have turned into irreplaceable assets to study
43 the spatial distribution and redistribution of species under given environmental circumstances,
44 and changes thereof (Elith & Leathwick, 2009). This broad group of statistical models relates
45 known species occurrences (or presence-absence) with information about the environmental
46 conditions at these locations (Guisan & Thuiller, 2005, Elith & Leathwick, 2009, Jiménez-Valverde
47 *et al.*, 2011). Based on statistically or theoretically derived response curves (Guisan *et al.*, 2017),
48 they aim to define the environmental niche and related geographical range in which organisms
49 can operate. Applications range from studies on the effects of anthropogenic climate change to
50 predictions of biological invasions, and they have proven valuable across a wide range of
51 organism groups both above and below the soil surface (Martiny *et al.*, 2006, Elith & Leathwick,
52 2009, Sato *et al.*, 2012, Bahram *et al.*, 2018). For belowground organisms in particular, however,
53 the successful application of HSDMs needs to overcome important limitations concerning the
54 scale, resolution and accuracy that can reliably and routinely be obtained. Indeed, many of the
55 available global, high-resolution gridded environmental variables cover the above-ground world

56 only, with many of the key drivers of soil microbial distributions (like soil pH and soil organic
57 matter) insufficiently recorded, or at coarse resolutions only (Fierer & Jackson, 2006, Mod *et al.*,
58 2016). This spatial mismatch between the environment as experienced by organisms on the one
59 hand, and the data we have available to model their distribution on the other, is particularly
60 pronounced for soil microbes. Their distribution, more than for most aboveground species groups,
61 indeed depends largely on processes happening at a small sub-meter scale in the hard-to-
62 measure belowground environment.

63

64 To understand and predict soil functioning in a rapidly changing world, we need the ability to
65 accurately model the spatial distribution patterns of soil microorganisms. To achieve this, we need
66 to incorporate local-scale drivers and processes relevant to these species groups in our HSDMs,
67 in order to improve predictions of large-scale distributional patterns, especially in response to
68 global change. For example, accurately attributing local variance to abiotic drivers versus biotic
69 interactions can aid our understanding of the importance of species interactions involving
70 microbes as drivers of spatial distributions (de Mesquita *et al.*, 2016, Ovaskainen *et al.*, 2017).
71 Such an approach can help improve our predictions of the faith of ecologically important, but so
72 far elusive groups including for example plant pathogens, mycorrhizal fungi and nitrogen fixing
73 bacteria and their effects on ecosystem processes (Martiny *et al.*, 2006). While many of these
74 questions have been covered in microbial distribution studies before, a clear framework is still
75 lacking to incorporate the relevant processes at the local scale of sub-centimeters to meters – in
76 which microbes operate – with the large-scale spatial patterns across several kilometers that are
77 the key focus in spatial ecology, without the need for an ever-increasing higher resolution of the
78 used spatial data.

79

80 **Belowground microbial distributions: the importance of the local scale**

81

82 It has long been assumed that many soil microbiota have cosmopolitan distributions. This theory
83 stems from the observation that they have large population sizes and short generation times
84 resulting in high dispersal, as well as the capacity to disperse over long distances if they get into
85 the air or water streams (Frey, 2015). Additionally, the actual drivers of their occurrence might
86 vary more on a sub-centimeter scale than on the kilometer-scale traditionally used, as is the case
87 for e.g. pore size, or biotic interactions with other soil microbiota. If soil microbial distributions
88 would indeed relate only little to global-level environmental gradients, the validity of HSDM-
89 approaches would be hampered, as it would be impossible to identify an environmental niche at

90 the regional or global scale. However, recent research has demonstrated that spatial patterns of
91 microbial diversity qualitatively similar to those observed for plants and animals exist, suggesting
92 that soil microbes indeed "have a biogeography" (Martiny *et al.*, 2006, Hanson *et al.*, 2012, Sato
93 *et al.*, 2012, Peay *et al.*, 2016). While the existence of such spatial structuring should not come
94 as a surprise for host-dependent soil microbes like plant-associated mycorrhizae, the presence
95 of a spatial distributional pattern is more telling for microbes like saprotrophs that operate
96 independent of a plant host. Nevertheless, their high dispersal flexibility implies that – even more
97 so than for most aboveground organisms - the range of environmental conditions and the nature
98 of biotic interactions at the local scale is likely to be of particular importance for their distribution
99 (Sato *et al.*, 2012, Vos *et al.*, 2013), and an efficient way to include this local environmental
100 variation in HSDMs is thus crucial.

101
102 While existing HSDMs often rely heavily on air temperature averages obtained from global climate
103 models (e.g. Aguilar & Lado, 2012, Tedersoo *et al.*, 2014), the use of disconnected environmental
104 variables overlooking the small-scale nature of the community processes involved is unlikely to
105 provide meaningful descriptions of the regional and global distribution of soil biota, and thus not
106 fully exploit the potential of HSDMs for soil organisms (Fierer & Jackson, 2006). Recent HSDMs
107 studies indeed increasingly acknowledge the need for more environmental data at the scale and
108 especially the location relevant to the studied organisms, as opposed to the traditional use of
109 coarse spatiotemporal averages of environmental conditions (Lembrechts *et al.*, 2018,
110 Lembrechts *et al.*, 2019, Zellweger *et al.*, 2019). This issue becomes particularly important when
111 one aims to use HSDMs predictively to investigate the effects of global (or local) change on
112 species (re)distributions. Incorporating local conditions is for example critical in regard to
113 microrefugia, where isolated populations can persist in a favorable microclimate amid
114 deteriorating climatic conditions until the latter become favorable again (Hannah *et al.*, 2014,
115 Lenoir *et al.*, 2017). For many soil microbiota, these microrefugia might exist within centimeters
116 to meters from where the species currently occurs due to the large local heterogeneity in
117 conditions (Veresoglou *et al.*, 2015), and HSDMs averaging out these local gradients are thus
118 likely to largely overestimate local extinction incidences, as has been shown for other organism
119 groups (Lenoir *et al.*, 2013). Similarly, this local heterogeneity can create stepping stones of
120 favorable conditions that facilitate range expansion beyond that predicted at a coarse scale
121 (Lembrechts *et al.*, 2017). The relevant data to counter these biased estimates is nevertheless
122 rarely available.

123

124 **The way forward to improve HSDMs for soil microbiota**

125

126 In this paper, we introduce a conceptual framework that provides the missing link between the
127 local-scale (abiotic and biotic) belowground environment and large-scale distribution patterns,
128 building on the recent advancements in tools, measurement techniques and models (Fig. 1, see
129 also Lembrechts (2020) for a simple overview of the followed steps in R). Keeping the strengths
130 and limitations of HSDMs in mind, this framework builds on 1) smart *in-situ* measurements of local
131 environmental heterogeneity and community distributions, 2) recent developments in spatial
132 modelling of environmental conditions, and 3) the hierarchical integration of the local with the
133 regional scale. It thus aims for a relatively simple way to incorporate the complexity of the
134 belowground world into statistical HSDMs without dramatically increasing the sampling effort, and
135 is especially suited for distributions driven by a complex interaction between local environmental
136 conditions and biotic interactions which mechanistic distribution models cannot (yet) handle.

137 Before expanding further, it is noteworthy that the definition of a 'species', as traditionally used in
138 HSDMs, is far less straightforward for soil microbes than it is for aboveground organisms, and
139 that approaches to identify micro-organism identities differ for different microbial groups. The
140 following framework focuses specifically on soil micro-organisms that are routinely identified as
141 operational taxonomic units (OTUs) or amplicon sequence variants (ASVs). These OTUs are
142 commonly used as proxies for species (e.g. Bahram *et al.*, 2016), yet recently, delineation of taxa
143 has often been abandoned in favor of ASVs (Knight *et al.*, 2018) as the unit of analysis, in
144 particular when sequence heterogeneity is not too high. The latter avoids the arbitrary binning of
145 sequences and allows comparisons across different studies, yet it would increase complexity of
146 the models due to a higher number of identities many of which would be intraspecific, as would
147 be the case for e.g. fungal ITS-sequences (Heeger *et al.*, 2018). If needed, meaningful HSDMs
148 can also still be produced by focusing on keystone taxa or higher taxonomic resolutions only.

149

150 ***Nested sampling of abiotic conditions***

151

152 Global databases and gridded datasets are currently being developed for many soil variables
153 (e.g. SoilGrids, Hengl *et al.*, 2017). Even though in most cases database resolutions are still
154 relatively coarse for application to the world of microbiota, the accuracy of the available data on
155 many local environmental drivers is rapidly rising. For soil moisture, for instance, the first attempts
156 to obtain global remotely sensed data at a coarse scale (3 and 9 km) are now appearing (He *et al.*,
157 *et al.*, 2015). For soil microbial distributions, however, large-scale gridded datasets are still several

158 orders of magnitude too coarse in spatial resolution. For many other relevant variables, e.g. soil
159 climate, users often still have to rely on the above-ground equivalent for higher resolutions (e.g.
160 free-air temperatures from CHELSA at 1 km² resolution (Karger *et al.*, 2017) vs soil temperatures
161 from ERA5 at 9 km² resolution (Copernicus Climate Change Service, 2019)).

162
163 The success of any application of HSDMs to soil microbial distributions thus fundamentally needs
164 to rely on a sampling design that allows linking the local heterogeneity, inherently present in soils,
165 with the microbial distribution at regional scales (Fig. 2). We therefore recommend a nested
166 sampling approach, measuring the relevant environmental variables right where it matters for
167 each specific study organism (e.g. measuring soil temperatures *in-situ* instead of relying on free-
168 air large-scale interpolations (Lembrechts *et al.*, 2018). The idea of a nested sampling design is
169 to capture the spatial autocorrelation signal that may exist both at a very fine spatial resolution as
170 well as across larger spatial extents, due to the modifying impact of both climatic and biotic
171 interactions on soil conditions (King *et al.*, 2008, King *et al.*, 2010, Ovaskainen *et al.*, 2017). Such
172 nested sampling can not only provide a better mechanistic understanding of drivers of the local
173 and regional distribution of a focal study organism, it is also a crucial requirement to use HSDMs
174 for predictions under global change (Mateo *et al.*, 2019).

175
176 These measurements should cover both the local (e.g. meters or even sub-centimeters apart, or
177 at different depths in the soil) and regional scales (e.g. up to kilometers apart) (King *et al.*, 2010).
178 Importantly, one should not aim for the finest possible resolution of environmental variables
179 across the whole landscape. Such a run for an increasing refinement in the spatial resolution
180 would greatly reduce the efficiency of HSDMs over larger spatial extents, while not necessarily
181 improving their accuracy (Bennie *et al.*, 2014, Lembrechts *et al.*, 2018). Instead, we suggest to
182 explicitly keep local-scale heterogeneity as a variable in regional-scale distribution models when
183 aggregating to a coarser resolution. This can for example be achieved by using environmental
184 data with a fine spatiotemporal resolution from a selection of sites to estimate the variation around
185 the mean within a certain measurement location both in space and time (Fig. 2a). For example,
186 one can take into account the differences in the variability in soil moisture levels occurring within
187 different environments along large-scale environmental gradients (e.g in mountains vs. flatlands).
188 Such estimates of within-pixel heterogeneity can be obtained with relatively simple aggregating
189 techniques, for example by using the local standard deviation in addition to the mean, or by fitting
190 correlative models that include within-pixel values as a covariate to incorporate uncertainty. If
191 desired, more complex approaches can be used, for example using mechanistic models that

192 describe the spatiotemporal variability in abiotic conditions in the system, or by calibrating HSDMs
193 at the local scale using Hierarchical Niche Models (HNMs) – distribution models that allow formally
194 including information from different spatial scales (Mateo *et al.*, 2019).

195

196 ***Converting nested sampling into gridded datasets***

197

198 Importantly, measured variables at the sampling locations are only one part of the story. To allow
199 spatial interpolations of a species' distribution between the sampling locations, gridded
200 environmental datasets – both of the means and the variation around these means – are
201 necessary. The absence of such gridded datasets for belowground environmental variables with
202 sufficient resolution and reliability is one of the fundamental issues currently hampering further
203 improvements in HSDMs, especially so for belowground organisms. In what follows, we describe
204 how the measured environmental conditions can be converted into gridded products with a spatial
205 resolution far finer than what is currently used, including the even finer-scaled local heterogeneity
206 at the sub-centimeter to meter scale.

207

208 Recent work has shown that remotely sensed drivers and proxies can accurately be linked as
209 covariates to measured environmental variables in the soil by using hybrid models combining
210 statistical correlations with mechanistic knowledge of the drivers of environmental variation (Fig.
211 2, Lembrechts *et al.*, 2018, Zellweger *et al.*, 2019). For example, soil moisture correlates strongly
212 with local topography, for which highly accurate data at meter to centimeter-resolution can be
213 obtained through remote sensing (satellite- or LiDAR-based DEMs) (Sørensen *et al.*, 2006). By
214 combining such remotely-sensed data with *in-situ* measurements, soil moisture can be modelled
215 with increasingly high accuracy, e.g. up to 1 m² resolution across a 3 km² regional extent
216 (Kemppinen *et al.*, 2017). Often, a resolution of 20 x 20 m up to 1 x 1 km is achievable using freely
217 available satellite-based remotely-sensed datasets, like digital elevation models. If one applies
218 these interpolation techniques to both the local averages and the within-pixel heterogeneity, both
219 can be converted into gridded datasets with the desired spatial resolution (Fig. 2b). To accomplish
220 this interpolation of the means and variation in these in-situ measurements, one can either use
221 simple statistical modelling methods - making sure that model residuals are not spatially
222 autocorrelated - such as generalized linear mixed-effects models, generalized additive mixed-
223 effects models, boosted regression trees, or random forests (e.g. Kemppinen *et al.*, 2017), as well
224 as spatially explicit geostatistical approaches such as spatial kriging or geographically weighted
225 regressions (GWR) (Fotheringham *et al.*, 2003, Lembrechts *et al.*, 2019). The latter technique

226 estimate model parameters at each geographical location by implementing a kernel and weighing
227 explanatory variables by distance. Although such geostatistical approaches have a more explicit
228 way of incorporating spatial structure than traditional generalized linear models or generalized
229 linear mixed-effects models for spatial interpolation of environmental variables, they cannot be
230 used to extrapolate microclimate outside the spatiotemporal extent covered by the data
231 (Lembrechts *et al.*, 2018). While doing this at the scale of single regional studies has thus proven
232 productive, globally coordinated efforts measuring and compiling such data on belowground
233 conditions – preferably using a nested framework - are needed more urgently than ever (Robock
234 *et al.*, 2000, Slessarev *et al.*, 2016, Hengl *et al.*, 2017). Such efforts are indeed required to make
235 a wider range of soil environmental variables available for interpolation at large spatial scales,
236 and to calibrate emerging mechanistic models that allow predicting variables at every relevant
237 scale (Kearney *et al.*, 2014).

238
239 Using these techniques shown to be effective in studies on aboveground organisms, in
240 combination with a nested sampling approach can help solve the persisting mismatch between
241 the available gridded environmental data products and the environment as it is perceived by soil
242 organisms themselves, and will greatly enhance the reliability of soil microbial HSDMs. It would
243 indeed allow to significantly improve predictive distribution models, as it permits for microrefugia
244 to be identified within a local pixel through the variation in observed values (Hattab *et al.*, 2014).
245 The number of samples needed to achieve this heavily depends on the local heterogeneity
246 present in the studied variable (e.g. a higher sampling density might be needed for soil moisture
247 in mountainous environments than in flat terrain).

248
249 Even if we obtain sufficiently detailed environmental variables to allow the calibration of accurate
250 and fine-grained HSDMs, the question remains which of these environmental variables are driving
251 the regional distribution of soil microbes. Importantly, which variables are relevant for a particular
252 group of microbes should be decided on a case-by-case basis, and requires a thorough
253 understanding of the microbial group under study. Indeed, these drivers – and the scale on which
254 they operate – will depend among others on the size of the studied organism groups. For example,
255 Birkhofer *et al.* (2012) showed that abiotic soil properties explained significant amounts of
256 variation in fungal diversity, but not in yeasts or bacteria. The use of HSDMs should thus always
257 be seen in parallel with other experimental approaches in which the mechanisms behind the
258 spatial drivers are further disentangled (Ettema & Wardle, 2002).

259

260 Nevertheless, we argue here that many available datasets of soil microbial diversity would be
261 suitable for our approach. The idea of nested sampling *per se* is not new in HSDMs (Diez &
262 Pulliam, 2007, Elith & Leathwick, 2009), and applications of similar sampling designs have been
263 used in studies of soil microbial diversity (i.e. repeated local sampling along environmental
264 gradients, or hierarchically structured designs to capture variation at different scales) (e.g. King
265 *et al.*, 2008, Bahram *et al.*, 2016, Zhou *et al.*, 2016). The current scale of most sampling
266 campaigns (collecting soil diversity and environmental information with an accuracy of meters)
267 would indeed already provide sufficient detail to answer a myriad of questions regarding the true
268 relationship between soil microbial spatial distributions and the local environment, especially if
269 one would refrain from the common practice to pool all samples out of a plot together to create
270 an ‘average’ sample, and thus allow objective integration of local-scale heterogeneity into models
271 (Manter *et al.*, 2010). Global database efforts compiling data on soil microbial distributions and/or
272 their abiotic drivers can for example explicitly keep this hierarchy in their structure and, as much
273 as possible, link up the diversity data with in-situ measured local-scale environmental conditions.
274 Emerging efforts to synchronize such sampling efforts across the globe (e.g. Kao *et al.*, 2012,
275 Thompson *et al.*, 2017) are thus pivotal to move this forward, as are efforts to standardize
276 sampling protocols between studies (Halbritter *et al.*, 2020) to reduce uncertainties when merging
277 studies (Ramirez *et al.*, 2018).

278

279 ***Pooling soil microbial communities***

280

281 For the sampling of soil microbial communities, practical restrictions (time and money) often
282 preclude the analysis of large numbers of samples. However, as summarized in Fig. 2c, repeated
283 sampling of local communities within a plot is not a prerequisite for the successful application of
284 our framework. Indeed, soil microbial studies traditionally measure the pooled community
285 structure by analyzing a sample of mixed soil from a random set of locations within a plot and
286 expressing OTUs by their (normalized) read number after sequencing (Staddon *et al.*, 1997,
287 Cleary *et al.*, 2016). This approach on its own already provides an integration of the variability in
288 local communities driven by the local environmental heterogeneity. Linking the pooled community
289 data – which is thus weighed based on the local distribution of environmental conditions - to the
290 average and heterogeneity in local environmental conditions effectively takes into account the
291 necessary fine-scale variation. It is however recommended to validate the local relationship
292 between microbial communities as displayed in Fig. 2b and local environmental conditions in a
293 few sites.

294

295 Importantly, note that OTU abundances are semi-quantitative relative abundances only. They are
296 indeed only meaningful in relation to the community in which they are measured, but they do give
297 information on shifts in dominance of a certain OTU between different samples. Additionally, the
298 proposed approach still suffers from false absences in the community data, as do all spatial
299 sampling schemes balancing output with resource input (Jiménez-Valverde & Lobo, 2006). The
300 latter issue is not solved here, and would not necessarily benefit proportionally from increasing
301 sample sizes.

302

303 ***Including biotic interactions in HSDMs of soil microbiota***

304

305 A recurring pattern in analyses of the spatial distribution of soil microbes is that the local spatial
306 distribution of many soil microbiota is not only determined by the abiotic environment, but also by
307 the presence or absence of other taxa (Ettema & Wardle, 2002, Wardle *et al.*, 2004, de Vries *et*
308 *al.*, 2012, Bahram *et al.*, 2018). Indeed, for many soil microbes spatial distribution largely depends
309 on where competitors, facilitators, mutualists, predators and/or other interactors reside, especially
310 at the local scale. This is for example the case for mycorrhizal fungi, which rely on the presence
311 of relevant host plants (and vice versa), but also strongly interact with each other, other fungi,
312 bacteria, and protists (Sato *et al.*, 2012). To get reliable predictions of their spatial distributions, it
313 will thus be important to include these local biotic interactions, within and between trophic levels,
314 in our HSDMs.

315 Again, there are a variety of ways in which one can integrate these community effects. One could
316 for example include plot-level metrics of the cohesion of interaction networks (Herren & McMahon,
317 2017) or metrics for food web complexity (Pellissier *et al.*, 2013) as covariates in HSDMs.
318 Alternatively, a recently developed species distribution modelling tool now also allows to integrate
319 the need for both an environmental and community-based approach in our HSDM-framework:
320 joint species distribution modelling (jSDMs, see e.g. Ovaskainen *et al.*, 2017, Tikhonov *et al.*,
321 2017, Fig. 1). The so-called jSDMs are a class of models that jointly analyze multiple species'
322 distributions, quantifying both species-specific environmental responses (i.e. the environmental
323 niche of the focal species) and covariance among species (i.e. positive or negative co-
324 occurrences that can potentially capture biotic interactions). The use of jSDMs has proven a major
325 step forward regarding the modelling of the distribution of a wide range of aboveground species
326 groups, especially when incorporating a hierarchical study design (Ovaskainen *et al.*, 2017). It is
327 important to remember, however, that jSDMs are not suitable for highly diverse communities of

328 more than 100 species (with example studies commonly using up to around 50 co-occurring
329 species, Ovaskainen *et al.*, 2017, Tikhonov *et al.*, 2017). One should thus focus on specific groups
330 of soil microbes, like nitrogen fixers, mycorrhizae or pathogens, and/or restrict analyses to the
331 most common OTUs only. The latter is realistic, as the structure of belowground communities is
332 often very uneven, with the top 2% of microbial taxa worldwide shown to make up 41% of microbial
333 abundance (Delgado-Baquerizo *et al.*, 2018). If one does want to model whole communities of
334 thousands of taxa, other network based approaches, e.g. using plot-level interaction network
335 metrics like the cohesion factor of interaction networks as mentioned above can be
336 recommended.

337

338 For optimization of HSDMs for soil microbiota, we strongly recommend this integration of both
339 abiotic and biotic drivers. Interestingly, this can work both for the interactions between soil
340 organisms among each other, and between soil organisms and the aboveground world (e.g.
341 between plant species and their mycorrhiza, if such data is available, see also Fig. 1). Vice versa,
342 jSDMs can also be used to include the role of soil microbes as drivers of the distribution of
343 aboveground organisms like plants (de Mesquita *et al.*, 2016). The result is a hybrid model that
344 allows giving equal importance to specific species and local environmental conditions as drivers
345 of species distributions.

346

347 Such integration of biotic interactions could also help bypass the issue of functional redundancy
348 in soil microbial ecology, i.e. the idea that many soil microbial taxa would share the same function
349 and thus are often only absent (or rare) because their redundant counterpart is common (Allison
350 & Martiny, 2008, Mori *et al.*, 2016). This would result in an apparent unfilling of their ecological
351 niche, i.e. the absence of the species in an environment with suitable environmental conditions.
352 However, if the presence of a certain redundant species affects the distribution of the focal
353 species, this will be picked up by the species correlation matrix used in the jSDM and result in a
354 high relative importance of biotic co-occurrences.

355

356 **Conclusion**

357

358 In this perspective paper, we argued that we should aim to model the current and projected
359 distribution of soil microbiota at a relevant scale, resolution and measurement location, and with
360 the relevant variables that drive their distributions. We propose using a nested sampling and
361 hierarchical modelling approach for the spatial niche and associated distribution, with the

362 possibility of including different drivers of the distribution of the studied microbiota at different
363 scales, and accounting for (above- and belowground) biotic interactions, and refer to the recent
364 advancements in high-resolution spatial environmental modelling. Applying this framework to
365 existing datasets of soil microbial distributions - and keeping it in mind when setting up sampling
366 networks and globally coordinated efforts - promises to resolve many remaining questions about
367 the local, regional and global distribution of soil microbes. Eventually this approach will allow
368 answering fundamental questions concerning microbial distribution and community assembly, but
369 will also be relevant in applied work such as predicting abundance and resilience of soil-dwelling
370 pests or beneficial soil microbes in agriculture. Importantly, this framework paves the way towards
371 the prediction of changes in soil microbial distributions as a result of different global change
372 drivers, for which first the true relationships between organismal distributions and the environment
373 need to be unraveled.

374

375 **Funding**

376

377 This work was supported by a postdoctoral fellowship from the Research Foundation – Flanders (FWO) to
378 J.J.L.

379

380 **Acknowledgements**

381

382 We gratefully acknowledge the support from the Research Foundation - Flanders (FWO) to J.J.L. and the
383 comments from the 3 reviewers.

384

385 **References**

386

387 Aguilar M & Lado C (2012) Ecological niche models reveal the importance of climate variability
388 for the biogeography of protosteloid amoebae. *The ISME journal* **6**: 1506.

389 Allison SD & Martiny JB (2008) Resistance, resilience, and redundancy in microbial
390 communities. *Proceedings of the National Academy of Sciences* **105**: 11512-11519.

391 Bahram M, Kohout P, Anslan S, Harend H, Abarenkov K & Tedersoo L (2016) Stochastic
392 distribution of small soil eukaryotes resulting from high dispersal and drift in a local environment.
393 *The ISME journal* **10**: 885.

394 Bahram M, Hildebrand F, Forslund SK, Anderson JL, Soudzilovskaia NA, Bodegom PM,
395 Bengtsson-Palme J, Anslan S, Coelho LP & Harend H (2018) Structure and function of the
396 global topsoil microbiome. *Nature* **560**: 233.

397 Bennie J, Wilson RJ, Maclean IMD & Suggitt AJ (2014) Seeing the woods for the trees - when is
398 microclimate important in species distribution models? *Global Change Biology* **20**: 2699-2700.

399 Birkhofer K, Schöning I, Alt F, Herold N, Klärner B, Maraun M, Marhan S, Oelmann Y, Wubet T
400 & Yurkov A (2012) General relationships between abiotic soil properties and soil biota across
401 spatial scales and different land-use types. *PLoS One* **7**: e43292.

402 Cleary DW, Bishop AH, Zhang L, Topp E, Wellington EM & Gaze WH (2016) Long-term
403 antibiotic exposure in soil is associated with changes in microbial community structure and
404 prevalence of class 1 integrons. *FEMS microbiology ecology* **92**.
405 Copernicus Climate Change Service (2019) C3S ERA5-Land reanalysis. p.^pp.
406 de Mesquita CPB, King AJ, Schmidt SK, Farrer EC & Suding KN (2016) Incorporating biotic
407 factors in species distribution modeling: are interactions with soil microbes important?
408 *Ecography* **39**: 970-980.
409 de Vries FT, Manning P, Tallwin JR, Mortimer SR, Pilgrim ES, Harrison KA, Hobbs PJ, Quirk
410 H, Shipley B & Cornelissen JH (2012) Abiotic drivers and plant traits explain landscape-scale
411 patterns in soil microbial communities. *Ecology letters* **15**: 1230-1239.
412 Delgado-Baquerizo M, Oliverio AM, Brewer TE, Benavent-González A, Eldridge DJ, Bardgett
413 RD, Maestre FT, Singh BK & Fierer N (2018) A global atlas of the dominant bacteria found in
414 soil. *Science* **359**: 320-325.
415 Diez JM & Pulliam HR (2007) Hierarchical analysis of species distributions and abundance
416 across environmental gradients. *Ecology* **88**: 3144-3152.
417 Elith J & Leathwick JR (2009) Species Distribution Models: ecological explanation and
418 prediction across space and time. *Annual Review of Ecology Evolution and Systematics*, Vol. 40
419 p.^pp. 677-697.
420 Ettema CH & Wardle DA (2002) Spatial soil ecology. *Trends in ecology & evolution* **17**: 177-183.
421 Fierer N & Jackson RB (2006) The diversity and biogeography of soil bacterial communities.
422 *Proceedings of the National Academy of Sciences* **103**: 626-631.
423 Fotheringham A, Brunsdon C & Charlton M (2003) *Geographically weighted regression: the*
424 *analysis of spatially varying relationships*. John Wiley & Sons, Hoboken, USA.
425 Frey SD (2015) The spatial distribution of soil biota. *Soil Microbiology, Ecology, and*
426 *Biochemistry*, p.^pp. 223-244. Academic Press London, UK.
427 Guisan A & Thuiller W (2005) Predicting species distribution: offering more than simple habitat
428 models. *Ecology Letters* **8**: 993-1009.
429 Guisan A, Thuiller W & Zimmermann NE (2017) *Habitat suitability and distribution models: with*
430 *applications in R*. Cambridge University Press.
431 Halbritter AH, De Boeck HJ, Eycott AE, Reinsch S, Robinson DA, Vicca S, Berauer B,
432 Christiansen CT, Estiarte M & Grünzweig JM (2020) The handbook for standardized field and
433 laboratory measurements in terrestrial climate change experiments and observational studies
434 (ClimEx). *Methods in Ecology and Evolution* **11**: 22-37.
435 Hannah L, Flint L, Syphard AD, Moritz MA, Buckley LB & McCullough IM (2014) Fine-grain
436 modeling of species' response to climate change: holdouts, stepping-stones, and microrefugia.
437 *Trends in Ecology & Evolution* **29**: 390-397.
438 Hanson CA, Fuhrman JA, Horner-Devine MC & Martiny JB (2012) Beyond biogeographic
439 patterns: processes shaping the microbial landscape. *Nature Reviews Microbiology* **10**: 497.
440 Hattab T, Albouy C, Lasram FBR, Somot S, Le Loc'h F & Leprieur F (2014) Towards a better
441 understanding of potential impacts of climate change on marine species distribution: a
442 multiscale modelling approach. *Global ecology and biogeography* **23**: 1417-1429.
443 He KS, Bradley BA, Cord Af, Rocchini D, Tuanmu MN, Schmidlein S, Turner W, Wegmann M &
444 Pettorelli N (2015) Will remote sensing shape the next generation of species distribution
445 models. *Remote Sensing in Ecology and Conservation* **1**: 4-18.
446 Heeger F, Bourne EC, Baschien C, Yurkov A, Bunk B, Spröer C, Overmann J, Mazzoni CJ &
447 Monaghan MT (2018) Long-read DNA metabarcoding of ribosomal RNA in the analysis of fungi
448 from aquatic environments. *Molecular ecology resources* **18**: 1500-1514.
449 Hengl T, de Jesus JM, Heuvelink GB, Gonzalez MR, Kilibarda M, Blagotić A, Shangguan W,
450 Wright MN, Geng X & Bauer-Marschallinger B (2017) SoilGrids250m: Global gridded soil
451 information based on machine learning. *PLoS one* **12**: e0169748.

452 Herren CM & McMahon KD (2017) Cohesion: a method for quantifying the connectivity of
453 microbial communities. *The ISME journal* **11**: 2426.

454 Jiménez-Valverde A, Peterson AT, Soberon J, Overton JM, Aragon P & Lobo JM (2011) Use of
455 niche models in invasive species risk assessments. *Biological Invasions* **13**: 2785-2797.

456 Jiménez-Valverde A & Lobo J (2006) The ghost of unbalanced species distribution data in
457 geographical model predictions. *Diversity and Distributions* **12**: 521-524.

458 Kao RH, Gibson CM, Gallery RE, Meier CL, Barnett DT, Docherty KM, Blevins KK, Travers PD,
459 Azuaje E & Springer YP (2012) NEON terrestrial field observations: designing continental-scale,
460 standardized sampling. *Ecosphere* **3**: 1-17.

461 Karger DN, Conrad O, Böhrner J, Kawohl T, Kreft H, Soria-Auza RW, Zimmermann NE, Linder
462 HP & Kessler M (2017) Climatologies at high resolution for the earth's land surface areas.
463 *Scientific Data* **4**: 170122.

464 Kearney MR, Isaac AP & Porter WP (2014) microclim: Global estimates of hourly microclimate
465 based on long-term monthly climate averages. *Scientific data* **1**: 140006.

466 Kemppinen J, Niittynen P, Riihimäki H & Luoto M (2017) Modelling soil moisture in a high-
467 latitude landscape using LiDAR and soil data. *Earth Surface Processes and Landforms*.

468 King AJ, Meyer A & Schmidt SK (2008) High levels of microbial biomass and activity in
469 unvegetated tropical and temperate alpine soils. *Soil Biology and Biochemistry* **40**: 2605-2610.

470 King AJ, Freeman KR, McCormick KF, Lynch RC, Lozupone C, Knight R & Schmidt SK (2010)
471 Biogeography and habitat modelling of high-alpine bacteria. *Nat Commun* **1**: 53.

472 Knight R, Vrbanac A, Taylor BC, Aksenov A, Callewaert C, Debelius J, Gonzalez A, Kosciulek
473 T, McCall L-I & McDonald D (2018) Best practices for analysing microbiomes. *Nature Reviews*
474 *Microbiology* **16**: 410.

475 Lembrechts J, Nijs I & Lenoir J (2018) Incorporating microclimate into species distribution
476 models. *Ecography*.

477 Lembrechts JJ (2020) A framework to bridge scales in distribution modelling of soil microbiota –
478 a modelled example. p.^pp.

479 Lembrechts JJ, Lenoir J, Nuñez MA, Pauchard A, Geron C, Bussé G, Milbau A & Nijs I (2017)
480 Microclimate variability in alpine ecosystems as stepping stones for non-native plant
481 establishment above their current elevational limit. *Ecography* **40**: 001-009.

482 Lembrechts JJ, Lenoir J, Roth N, Hattab T, Milbau A, Haider S, Pellissier L, Pauchard A, Ratier
483 Backes A & Dimarco RD (2019) Comparing temperature data sources for use in species
484 distribution models: From in-situ logging to remote sensing. *Global Ecology and Biogeography*
485 **28**: 1578-1596.

486 Lenoir J, Hattab T & Pierre G (2017) Climatic microrefugia under anthropogenic climate change:
487 implications for species redistribution. *Ecography* **40**: 253-266.

488 Lenoir J, Graae BJ, Aarrestad PA, *et al.* (2013) Local temperatures inferred from plant
489 communities suggest strong spatial buffering of climate warming across Northern Europe.
490 *Global Change Biology* **19**: 1470-1481.

491 Manter DK, Weir TL & Vivanco JM (2010) Negative effects of sample pooling on PCR-based
492 estimates of soil microbial richness and community structure. *Appl Environ Microbiol* **76**: 2086-
493 2090.

494 Martiny JBH, Bohannan BJ, Brown JH, Colwell RK, Fuhrman JA, Green JL, Horner-Devine MC,
495 Kane M, Krumins JA & Kuske CR (2006) Microbial biogeography: putting microorganisms on
496 the map. *Nature Reviews Microbiology* **4**: 102.

497 Mateo RG, Aroca-Fernández MJ, Gastón A, Gómez-Rubio V, Saura S & García-Viñas JI (2019)
498 Looking for an optimal hierarchical approach for ecologically meaningful niche modelling.
499 *Ecological Modelling* **409**: 108735.

500 Mod HK, Scherrer D, Luoto M & Guisan A (2016) What we use is not what we know:
501 environmental predictors in plant distribution models. *Journal of Vegetation Science* **27**: 1308-
502 1322.

503 Mori AS, Isbell F, Fujii S, Makoto K, Matsuoka S & Osono T (2016) Low multifunctional
504 redundancy of soil fungal diversity at multiple scales. *Ecology letters* **19**: 249-259.

505 Ovaskainen O, Tikhonov G, Norberg A, Guillaume Blanchet F, Duan L, Dunson D, Roslin T &
506 Abrego N (2017) How to make more out of community data? A conceptual framework and its
507 implementation as models and software. *Ecology Letters* **20**: 561-576.

508 Peay KG, Kennedy PG & Talbot JM (2016) Dimensions of biodiversity in the Earth mycobiome.
509 *Nature Reviews Microbiology* **14**: 434.

510 Pellissier L, Rohr RP, Ndiribe C, Pradervand JN, Salamin N, Guisan A & Wisz M (2013)
511 Combining food web and species distribution models for improved community projections.
512 *Ecology and evolution* **3**: 4572-4583.

513 Ramirez KS, Knight CG, De Hollander M, Brearley FQ, Constantinides B, Cotton A, Creer S,
514 Crowther TW, Davison J & Delgado-Baquerizo M (2018) Detecting macroecological patterns in
515 bacterial communities across independent studies of global soils. *Nature microbiology* **3**: 189.

516 Robock A, Vinnikov KY, Srinivasan G, Entin JK, Hollinger SE, Speranskaya NA, Liu SX &
517 Namkhai A (2000) The Global Soil Moisture Data Bank. *Bulletin of the American Meteorological*
518 *Society* **81**: 1281-1299.

519 Sato H, Tsujino R, Kurita K, Yokoyama K & Agata K (2012) Modelling the global distribution of
520 fungal species: new insights into microbial cosmopolitanism. *Molecular Ecology* **21**: 5599-5612.

521 Slessarev E, Lin Y, Bingham N, Johnson J, Dai Y, Schimel J & Chadwick O (2016) Water
522 balance creates a threshold in soil pH at the global scale. *Nature* **540**: 567.

523 Sørensen R, Zinko U & Seibert J (2006) On the calculation of the topographic wetness index:
524 evaluation of different methods based on field observations. *Hydrology and Earth System*
525 *Sciences Discussions* **10**: 101-112.

526 Staddon W, Duchesne L & Trevors J (1997) Microbial diversity and community structure of
527 postdisturbance forest soils as determined by sole-carbon-source utilization patterns. *Microbial*
528 *Ecology* **34**: 125-130.

529 Tedersoo L, Bahram M, Pölme S, Kõljalg U, Yorou NS, Wijesundera R, Ruiz LV, Vasco-
530 Palacios AM, Thu PQ & Suija A (2014) Global diversity and geography of soil fungi. *science*
531 **346**: 1256688.

532 Thompson LR, Sanders JG, McDonald D, Amir A, Ladau J, Locey KJ, Prill RJ, Tripathi A,
533 Gibbons SM & Ackermann G (2017) A communal catalogue reveals Earth's multiscale microbial
534 diversity. *Nature* **551**: 457.

535 Tikhonov G, Abrego N, Dunson D & Ovaskainen O (2017) Using joint species distribution
536 models for evaluating how species-to-species associations depend on the environmental
537 context. *Methods in Ecology and Evolution* **8**: 443-452.

538 Veresoglou SD, Halley JM & Rillig MC (2015) Extinction risk of soil biota. *Nat Commun* **6**.

539 Vos M, Wolf AB, Jennings SJ & Kowalchuk GA (2013) Micro-scale determinants of bacterial
540 diversity in soil. *FEMS microbiology reviews* **37**: 936-954.

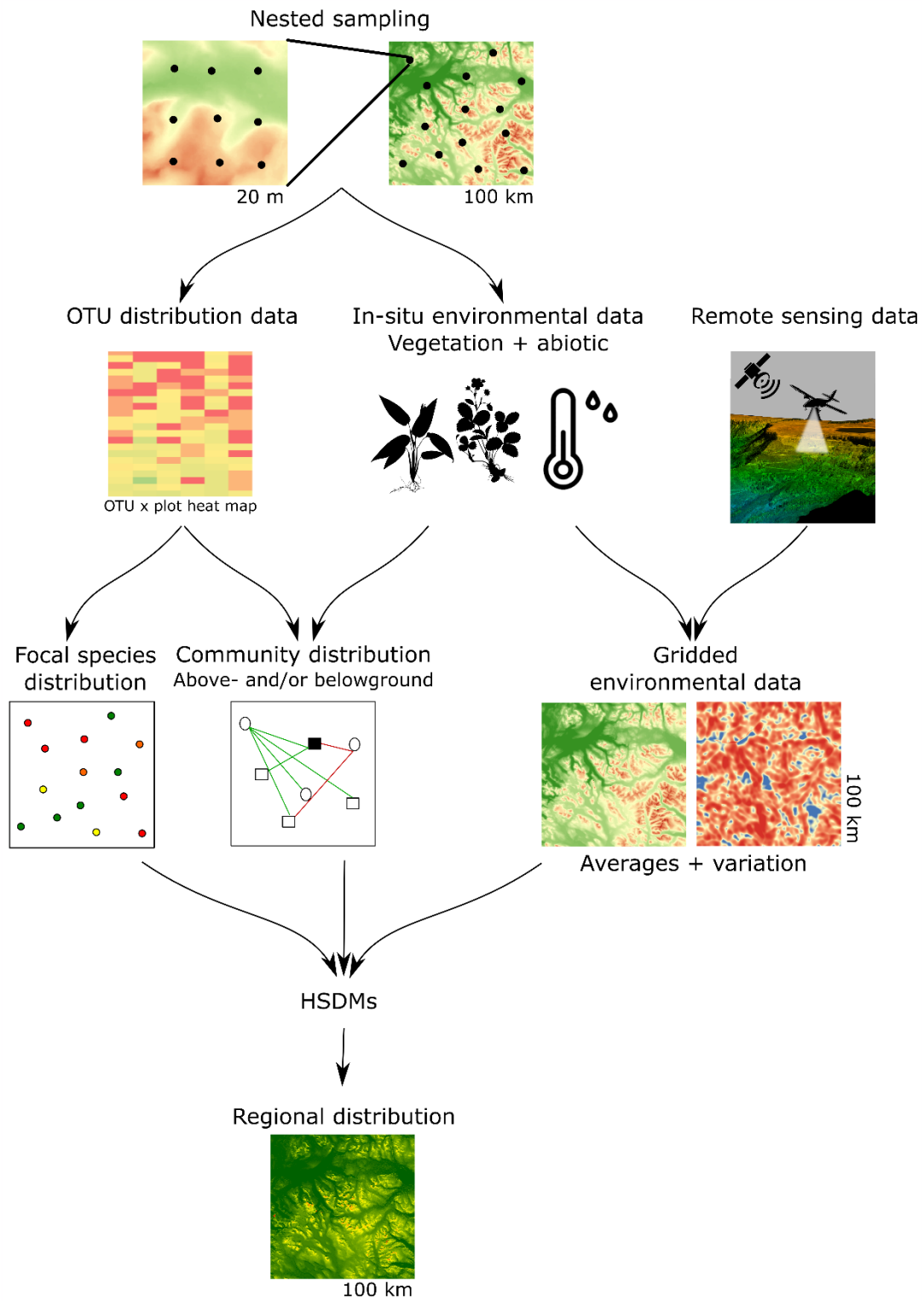
541 Wardle DA, Bardgett RD, Klironomos JN, Setälä H, Van Der Putten WH & Wall DH (2004)
542 Ecological linkages between aboveground and belowground biota. *Science* **304**: 1629-1633.

543 Zellweger F, De Frenne P, Lenoir J, Rocchini D & Coomes D (2019) Advances in microclimate
544 ecology arising from remote sensing. *Trends in ecology & evolution* **34**: 327-341.

545 Zhou J, Deng Y, Shen L, Wen C, Yan Q, Ning D, Qin Y, Xue K, Wu L & He Z (2016)
546 Temperature mediates continental-scale diversity of microbes in forest soils. *Nat Commun* **7**:
547 12083.

548

549



552 Fig. 1

553 *Figure 1. General overview of the proposed roadmap for HSDMs of soil microbiota at relevant spatial scales.*
554 *The strength of the approach lies in the nested sampling of above- and belowground communities and*
555 *environmental data, and the consequent inclusion of those at relevant scales in HSDMs. Microbial OTUs*
556 *and environmental conditions are sampled both at a high resolution in a selection of plots, and at a coarse*
557 *resolution across a region (see Fig. 2). The OTU-distribution data is then used to obtain information on the*
558 *distribution of a specific focal species of interest, as well as on its co-occurrences with other microbes in*
559 *the microbial community. The environmental variables can provide information on interactions with*
560 *aboveground organisms (e.g. between plants and their associated mycorrhizas), and on the species' abiotic*
561 *niche. The latter can be converted into gridded data, both for coarser-scale averages and their fine-scale*
562 *heterogeneity, by interpolating them using remotely sensed gridded environmental proxies. Then, the*
563 *regional distribution of specific focal species of interest can be modelled using habitat suitability and*
564 *distribution models (HSDMs), as a function of both these biotic interactions and the abiotic drivers, for*
565 *example using joint species distribution models (see main text). Note that displayed details regarding*
566 *sample sizes, resolution and spatial distribution of sample points are examples only, with actual*
567 *requirements depending on the study system, study question and practical limitations.*
568

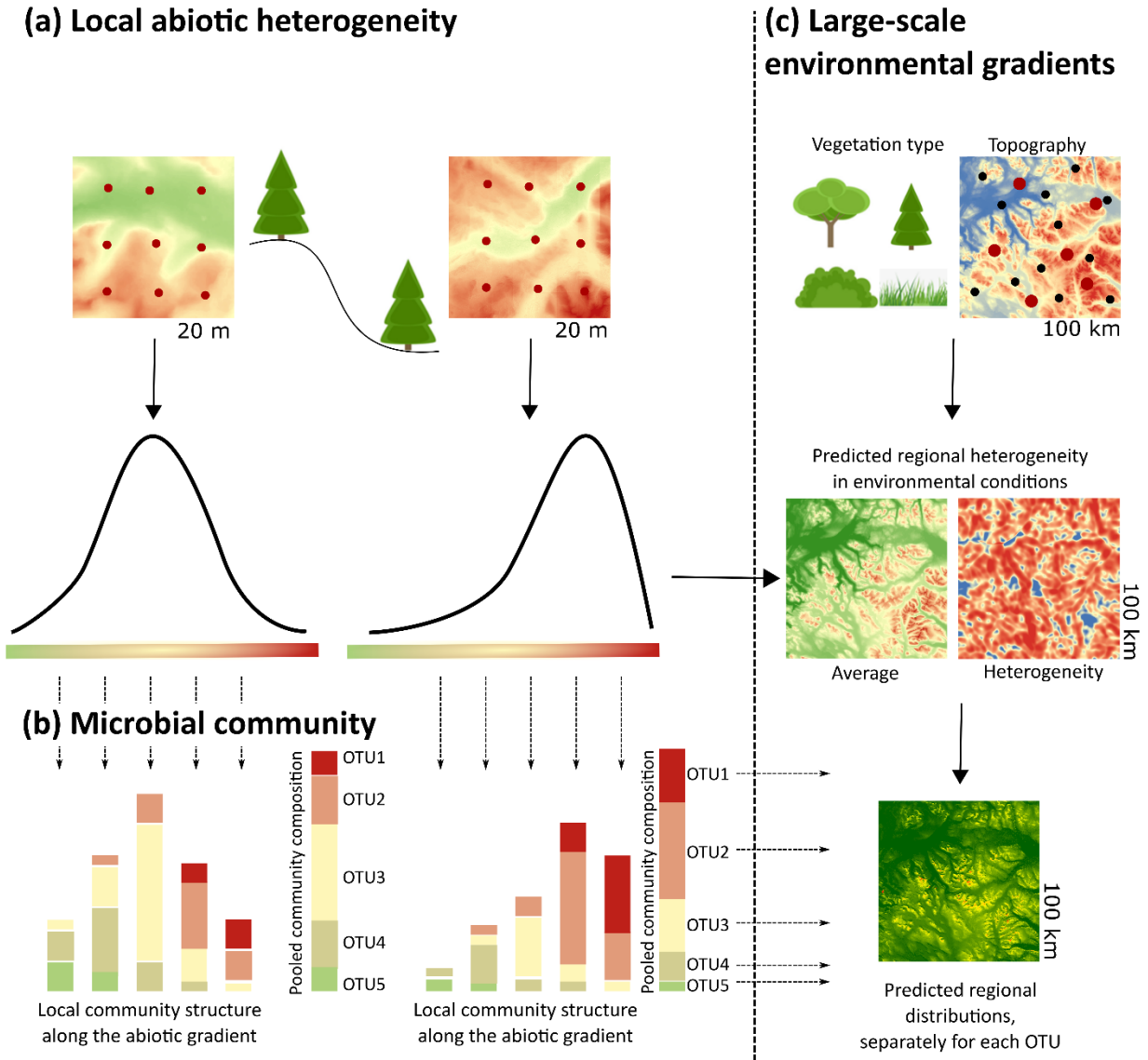


Fig. 2

569
570

571 *Figure 2. Detailed view of how the proposed nested sampling approach can be used to improve HSDMs,*
 572 *here exemplified for the link between microbial communities and local abiotic conditions (e.g. soil moisture).*
 573 *(a) Fine-scale measurements of the abiotic conditions of interest in a selection of small plots, distributed*
 574 *evenly along large-scale environmental gradients (e.g. in similar vegetation types, across topographic*
 575 *gradients in landscapes of 100 x 100 km), can be used to get an estimate of the local heterogeneity in*
 576 *environmental conditions. (b) These local relationships in a selection of plots (red dots) can be extrapolated*
 577 *across a whole region using coarser-grained averages of the same environmental variables (measured in*
 578 *the sites represented by black dots) and the relationships between these local conditions (averages and*
 579 *heterogeneity) and the large-scale environmental gradients. For example, one can extrapolate the*

580 *relationships in (a) using high-resolution digital elevation models and, if repeated in different vegetation*
581 *types, gridded vegetation maps. (c) The locally present microbial community depends on the local range of*
582 *abiotic conditions and the specific abiotic niche of each OTU (visualized here using the same color gradient*
583 *as for the environmental gradient). A pooled community composition sample will then have OTU*
584 *abundances dependent on the relative frequency of each local environmental condition within a plot. As*
585 *this pooled community data thus inherently contains the underlying local-scale variation in the microbial*
586 *community, one can calibrate distribution models using pooled community data measured in the red and*
587 *black dots depicted in (b). Finally, the distribution of each OTU can be modelled as a function of the*
588 *predicted regional averages and heterogeneity in abiotic conditions using HSDMs.*
589