

Genomic evidence for the wide-spread presence of lignocellulases among soil invertebrates

Hannah Muelbaier¹, Freya Arthen¹, Gemma Collins¹, Thomas Hickler¹, Karin Hohberg¹, Ricarda Lehmitz¹, Yannick Pauchet², Markus Pfenninger³, Anton Potapov¹, Julianne Romahn¹, Ina Schaefer¹, Stefan Scheu⁴, Clement Schneider⁵, Ingo Ebersberger⁶, and Miklós Bálint⁷

¹Affiliation not available

²Max-Planck-Institute for Chemical Ecology

³Senckenberg Biodiversität & Klima Forschungszentrum

⁴University of Goettingen

⁵Senckenberg Gesellschaft für Naturforschung

⁶Goethe University

⁷Senckenberg Biodiversity and Climate Research Centre

September 6, 2023

Abstract

Lignocellulose is a major component of plant biomass. Its decomposition is crucial for the terrestrial carbon cycle. Microorganisms are considered as primary decomposers and evidence increases that some invertebrates may also decompose lignocellulose. We investigated the taxonomic distribution and evolutionary origins of GH45 cellulases in a collection of soil invertebrate genomes and found that these genes are common in springtails and oribatid mites. Phylogenetic analysis revealed that cellulase genes were acquired early in the evolutionary history of these groups. Domain architectures and predicted 3D enzyme structures indicate that these cellulases are functional. Patterns of presence and absence of these genes across different lineages prompt further investigation into their evolutionary and ecological benefits. The ubiquity of cellulase genes suggests that soil invertebrates may play a role in lignocellulose decomposition, independently from microorganisms. Understanding the ecological and evolutionary implications might be crucial for understanding soil food webs and the carbon cycle.

1 **Genomic evidence for the wide-spread presence of**
2 **lignocellulases among soil invertebrates**

3

4 Hannah Muelbaier^{1,2}, Freya Arthen^{1,2}, Gemma Collins,^{2,3} Thomas Hickler^{4,5}, Karin Hohberg⁶,
5 Ricarda Lehmitz^{2,6}, Yannick Pauchet⁷, Markus Pfenninger^{2,4,8}, Anton Potapov^{9,10,11}, Juliane
6 Romahn^{2,4}, Ina Schaefer^{2,4,11}, Stefan Scheu¹², Clément Schneider^{2,6}, Ingo Ebersberger^{1,2,4}†,
7 Miklós Bálint^{2,4,13}*

8 † These authors contributed equally to this work

9 Affiliations:

10 ¹ Applied Bioinformatics Group, Inst. of Cell Biology and Neuroscience, Goethe University,
11 Max-von-Laue Str. 13, 60438 Frankfurt am Main, Germany

12 ² LOEWE Centre for Translational Biodiversity Genomics, Senckenbergenlage 25, 60325
13 Frankfurt am Main, Germany

14 ³ Manaaki Whenua – Landcare Research, Auckland, New Zealand

15 ⁴ Senckenberg Biodiversity and Climate Research Centre, Senckenbergenlage 25, 60325
16 Frankfurt am Main, Germany

17 ⁵ Department of Physical Geography, Goethe University, Altenhöferallee 1, 60438
18 Frankfurt/Main, Germany

19 ⁶ Senckenberg Museum of Natural History Görlitz, Am Museum 1, 02826 Görlitz, Germany

20 ⁷ Insect symbiosis, Max Planck Institute for Chemical Ecology, Beutenberg campus, Hans-
21 Knoell-Str. 8, 07745 Jena

22 ⁸ Institute for Molecular and Organismic Evolution, Johannes Gutenberg University, Mainz,
23 Germany

24 ⁹ German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig,
25 Puschstrasse 4, 04103, Leipzig, Germany

26 ¹⁰ Institute of Biology, University of Leipzig, Puschstrasse 4, 04103, Leipzig, Germany

27 ¹¹ Animal Ecology, University of Goettingen, Untere Karspüle 2, 37073, Goettingen,
28 Germany

29 ¹² "J.F. Blumenbach Institute of Zoology and Anthropology, University of Goettingen, Berliner
30 Str. 28 37073 Goettingen, Germany

31 ¹³ Institute of Insect Biotechnology, Justus-Liebig University, Heinrich-Buff-Ring 26, 35392,
32 Giessen, Germany

33 * miklos.balint@senckenberg.de

34

35 **Abstract**

36 Lignocellulose is a major component of plant biomass. Its decomposition is crucial for the
37 terrestrial carbon cycle. Microorganisms are considered as primary decomposers and
38 evidence increases that some invertebrates may also decompose lignocellulose. We
39 investigated the taxonomic distribution and evolutionary origins of GH45 cellulases in a
40 collection of soil invertebrate genomes and found that these genes are common in
41 springtails and oribatid mites. Phylogenetic analysis revealed that cellulase genes were
42 acquired early in the evolutionary history of these groups. Domain architectures and
43 predicted 3D enzyme structures indicate that these cellulases are functional. Patterns of
44 presence and absence of these genes across different lineages prompt further investigation
45 into their evolutionary and ecological benefits. The ubiquity of cellulase genes suggests that
46 soil invertebrates may play a role in lignocellulose decomposition, independently from
47 microorganisms. Understanding the ecological and evolutionary implications might be crucial
48 for understanding soil food webs and the carbon cycle.

49 **Introduction**

50 Most photosynthetically bound carbon on land ends up in woody plants as lignocellulose, a
51 composite of several polysaccharides. The decomposition of lignocellulose occurs
52 predominantly in soils, which returns most of this carbon into the atmosphere (Post et al.,
53 1990). Nevertheless, terrestrial ecosystems currently sequester about 29% of the
54 anthropogenic carbon emissions, which implies an important but not fully understood role of
55 terrestrial carbon cycling for climate regulation (Cragg et al., 2015). Microorganisms,
56 especially bacteria and fungi, encode glycoside hydrolase cocktails for lignocellulose
57 degradation in their genomes (Cragg et al., 2015), and are considered the main actors of
58 decomposition (Bradford et al., 2017; Crowther et al., 2019; Pausas and Bond, 2020). The
59 contribution of animals to decomposition of lignocellulose - beyond purely mechanical
60 shredding - remains less understood. Experiments have shown that the presence of soil
61 invertebrates can increase litter mass loss by up to 50 percent (García-Palacios et al., 2013).
62 It is estimated that they decompose more deadwood in tropical forests than free-living
63 microorganisms (Griffiths et al., 2019). Nevertheless, the mechanisms behind decomposition
64 performed by soil invertebrates remains obscure and the ability of soil animals to degrade
65 composite polysaccharides without relying on gut symbionts remains a long-standing debate
66 in soil ecology (Berg et al., 2004; Cragg et al., 2015).

67 It was originally assumed that lignocellulose degradation by animals was entirely
68 'outsourced' to the gut microbiome (Briones, 2018; García-Palacios et al., 2013). However,
69 evidence is emerging that at least some invertebrates, such as molluscs, crustaceans and
70 phytophagous insects can synthesize cellulase enzymes themselves (Busch et al., 2019;
71 Chang and Lai, 2018; Cragg et al., 2015; Griffiths et al., 2021; Han et al., 2022; Kern et al.,
72 2013; King et al., 2010; Shelomi et al., 2014; Watanabe et al., 1998). Scattered evidence
73 also exists for the expression of active endogenous cellulases by distantly-related soil
74 invertebrates, e.g. the earthworm *Pheretima hilgendorfi* (Nozaki et al., 2009), the Antarctic
75 springtail *Cryptopygus antarcticus* (Hong et al., 2014), as well as few oribatid mites and
76 other springtails (Busch et al., 2019). Based on these individual findings, we hypothesize
77 that a larger fraction of soil invertebrates than previously thought may be directly contributing
78 to lignocellulose decomposition of dead plant matter in soils, without necessarily relying on a

79 specialized microbiome. Given their global abundance and diversity in many soil ecosystems
80 (FAO et al., 2020; Phillips et al., 2019; Potapov et al., 2023; van den Hoogen et al., 2019),
81 soil invertebrates could therefore have an important but so far overlooked role in the
82 terrestrial carbon cycle which is distinct from the lignocellulose decomposition ability of
83 microorganisms. To evaluate whether endogenous decomposition ability is a common
84 feature shared by the main groups of soil invertebrates, we screened a diverse set of newly
85 sequenced genomes of Collembola, Enchytraeidae, Gamasina, Myriapoda, Nematoda,
86 Oribatida, Tardigrada as representatives of dominating and ubiquitous soil invertebrates, for
87 the presence and origin of cellulase genes.

88 Results

89 We used a fungal sequence as a starting point to obtain a comprehensive overview of the
90 taxonomic distribution of GH45-type cellulases across all species represented by a genome
91 in NCBI RefSeq. We identified orthologs to fungal GH45 cellulases in 16,401 bacteria, 910
92 archaea and 1,101 eukaryotes (Table S1). The resulting phylogenetic profile revealed that
93 GH45 cellulases are abundant only in fungi, where they are, however, limited to individual
94 systematic groups (Fig. 1A). In bacteria, only 31 species were found to harbour a GH45-type
95 cellulase. Next to fungi, orthologs to fungal cellulases were largely confined to animals (Fig.
96 1A). We found no evidence that these comprise fungal contaminations (see Taxonomic
97 assignment and contaminant detection in Materials and Methods). In metazoans, most
98 orthologs were found in arthropods, and they were completely absent in vertebrates.

99 The analysis of publicly available genomes shed light on the general distribution of GH45-
100 type cellulases across the Tree of Life. However, it lacked the resolution to investigate the
101 occurrence of this cellulase in soil invertebrates. We extended the analysis by including
102 novel genome assemblies for an additional 176 species representing a diverse selection of
103 soil invertebrates (NCBI BioProject PRJNA758215) into the analysis. This revealed a high
104 occurrence of GH45-type cellulases in springtails (70%; 56 out of 78 analyzed species) and
105 in oribatid mites (60%; 33 out of 54 species, Table S2). Additionally, we detected cellulases
106 in Coleoptera and Thysanoptera (2 out of 2 species, Table S1). In three out of nine
107 nematode species we also found GH45 cellulases, whereas none were found in 30
108 representatives of Chilopoda and Diplopoda (Table S2).

109 We found GH45-type cellulases in three of four known main springtail lineages
110 (Poduromorpha, Entomobryomorpha, Symphyleona), missing only from the earliest
111 branching Neelipleona. Within Symphyleona and Entomobryomorpha, cellulases are
112 consistently absent in one clade each (Supp. Fig S1).

113 GH45-type cellulases are present in almost all representatives of the basal clades in
114 Oribatida (Enarthronota, Mixonomata, Holosomata; Table S2). In contrast, they were missing
115 in half of the investigated species from the later-branching Brachypylina (Supp. Fig S2). All
116 of the species investigated here that did not have GH45 cellulase genes were sexually
117 reproducing, while the GH45 cellulase genes were present in other sexually-reproducing
118 species and parthenogens. Besides Oribatida, we investigated a second mite taxon, the
119 Gamasina. Gamasina is represented by two species in the present data set (Table S2) and
120 GH45 cellulase genes were absent in both.

121 We reconstructed the phylogenetic relationships of the sequences identified both in RefSeq
122 assemblies and in soil invertebrate genomes (Fig. 1B) to better understand the evolutionary
123 trajectory resulting in the present-day distribution of soil invertebrate GH45 cellulases.
124 Animals are paraphyletic in this tree, but most of the invertebrate GH45 orthologs (94 out of
125 107) were placed in only four phylogenetically largely homogeneous clades, one each for the
126 Collembola, Oribatida, Thysanoptera, and Coleoptera. On a larger scale, we found that the
127 invertebrate cellulase clades are embedded into an evolutionary background formed by
128 fungal and bacterial sequences. We noted that oribatids are placed in a single clade together
129 with a few interspersed bacterial sequences. We compared the phylogeny from oribatid
130 mites and springtails with our reconstructed gene tree (Fig. 1C). A reconstruction of GH45
131 phylogeny with all found cellulase co-orthologs shows a highly dynamic and complex
132 evolutionary history with lineage-specific gene duplications and losses (Supp. Fig.S3).

133 We investigated the GH45 domain architecture and 3D protein structure and found that
134 domain architecture in fungi, antarctic springtail and mustard leaf beetle were similar. The
135 enzymes in these species also had similar predicted 3D structures with the cellulases found
136 by us (Fig. 1D).

137 Discussion

138 Most cellulases discovered to date in metazoan genomes belong to the GH45 family (Busch
139 et al., 2019), endo- β -1,4,-glucanases (EC 3.2.1.4) which catalyze the decomposition of
140 complex cellulose into more accessible oligosaccharides (Davies et al., 1993). It was
141 hypothesized that invertebrate GH45 cellulases were repeatedly obtained via horizontal
142 gene transfer (HGT) from fungal donors (Busch et al., 2019). Our search of 18,412 RefSeq
143 genomes revealed that only 31 bacterial species harbour a GH45-type cellulase. This
144 confirms the earlier hypothesis about the evolutionary roots of this enzyme within fungi
145 (Busch et al., 2019). The taxonomic distribution of GH45 presence was non-uniform in
146 metazoans, with most orthologs being found in arthropods. Our results confirm previous
147 research which did not find evidence for GH45 presence in vertebrates (Chang and Lai,
148 2018).

149 The screening of new soil invertebrate genomes uncovered novel GH45 cellulase presence
150 patterns. GH45 genes were found in well over half of investigated springtail and oribatid mite
151 species. Springtails form a basal hexapod group abundant across the globe, especially in
152 cold regions (Potapov et al., 2023). They are known as fungal feeders, but also consume
153 detritus and fresh plant materials (Potapov et al., 2020). While the ortholog search showed
154 the presence of GH45-type cellulase in springtail lineages Poduromorpha,
155 Entomobryomorpha and Symphypleona it is missing from the earliest branching Neelipleona.
156 Neelipleona feed on detritus colonized by fungi and bacteria rather than on plant remains
157 (Potapov et al., 2022), and thus may not directly require lignocellulase activity. A single
158 acquisition event post-dating the divergence of the Neelipleona may explain this observation.
159 However, since Neelipleona are represented only by a few taxa, the conclusion on cellulase
160 absence should be taken cautiously. Within Symphypleona and Entomobryomorpha,
161 cellulases are consistently absent in one clade each. The absence of cellulase genes is
162 unexpected in the lucerne flea *Sminthurus viridis* that feeds on live leaf tissue (Greenslade
163 and Ireson, 1986), and in the families of Dicyrtomidae, Bourletiellidae and Sminthuridae
164 which consume mainly fresh plant materials (Potapov et al., 2020, 2022). The lack of

165 cellulase genes suggests that these taxa either outsource lignocellulose decomposition to
166 their midgut microbiome, rely on other GH families, or do not digest cellulose. However, the
167 presence of GH45-type cellulase genes seems to be a common trait among most springtails.

168 Mites are the most numerous arthropods on land (Rosenberg et al., 2023), with most
169 representatives in soil ecosystems belonging to Oribatida. Oribatid mites can have diverse
170 feeding strategies, but they mostly feed on leaf litter at different decomposition stages and
171 on microorganisms (Maraun et al., 2023). Our results showed that GH45-type cellulases are
172 present all over the basal clades of Oribatida while they were missing in half of the
173 Branchypylina, which are later-branching. One interesting finding is that GH45 cellulase is
174 missing especially in sexually reproducing mites, while they are present in other sexually
175 reproducing taxa. This fits to the pattern that parthenogenetic oribatid mites tend to occupy
176 lower trophic positions and typically function as primary decomposers, opposed to
177 secondary decomposers feeding predominantly on microorganisms (Fischer et al., 2014).
178 However, ecological interpretation of these patterns is difficult since we do not know if
179 species without GH45 cellulase genes contain other classes of cellulases, digest cellulose
180 with the help of their microbiome or are indeed incapable of cellulose digestion. In general,
181 the complex pattern of GH45 presence is similar to the low phylogenetic conservatism of
182 ecological traits in oribatids, such as feeding mode (Potapov et al., 2022).

183 As expected, we detected cellulases in Coleoptera (Kirsch et al., 2014). GH45-type
184 cellulases were completely absent in Chilopoda and Diplopoda. The latter was surprising as
185 Diplopoda are a key litter-feeding soil invertebrate group (Joly et al., 2020; Potapov et al.,
186 2022). GH45 cellulases were also absent in Gamasina mites which are predators and
187 therefore might not benefit from cellulose degradation. The first report of endogenous
188 cellulases in Thysanoptera suggests that our analysis uncovers only the tip of the iceberg.
189 We expect that taxonomically broad genome sequencing of eukaryotes promoted e.g. by the
190 Earth BioGenome Project (Formenti et al., 2022; Lewin et al., 2022) will recover further
191 animal groups in possession of enzymes targeting lignocellulose decomposition.

192 Taken together, our data suggests an early acquisition of a GH45-type cellulase during the
193 diversifications of both springtails and oribatids, instead of repeated horizontal transfer
194 events. This implies that the possession of a GH45-type cellulase is an ancestral trait in
195 these groups. Similar to our results, cellulase acquisition was shown to be important for the
196 diversification of herbivorous beetles (Kirsch et al., 2014). Differences in the GH45 cellulase
197 gene tree from the oribatid and springtail species trees likely result from a highly dynamic
198 evolution of the GH45 cellulase repertoire. Lineage-specific gene duplications and losses
199 have partially disconnected the evolutionary history of the contemporary cellulase genes
200 from the phylogeny of the species they are found in (Supp. Fig. S3). Lineage-specific
201 duplications have been described for other cellulases (Shelomi et al., 2016; Shin et al.,
202 2022), and differential duplicate loss has been shown to result in gene tree - species tree
203 incongruencies (Parey et al., 2020). The presence of cellulases detected in thrips suggests
204 that similar processes might have been important also during the evolution of other
205 arthropod groups.

206 Although cellulase presence in genomes is not a proof of function, several lines of evidence
207 point toward functionality. First, domain architecture of GH45 cellulases in fungi, antarctic
208 springtail and mustard leaf beetle are similar, and the enzymes themselves also have similar
209 predicted 3D structures with the cellulases found by us (Fig. 1D). Second, fungal (Cragg et

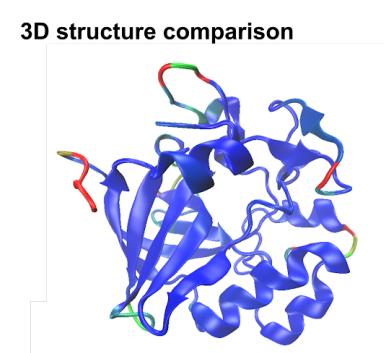
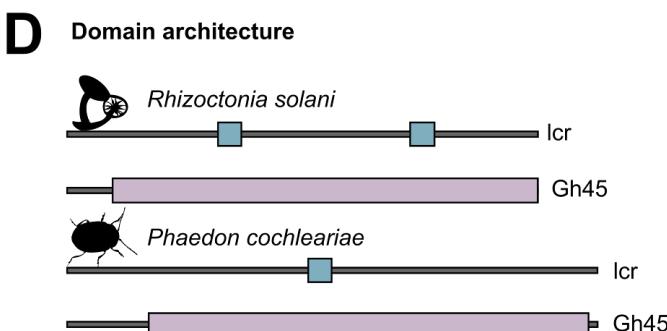
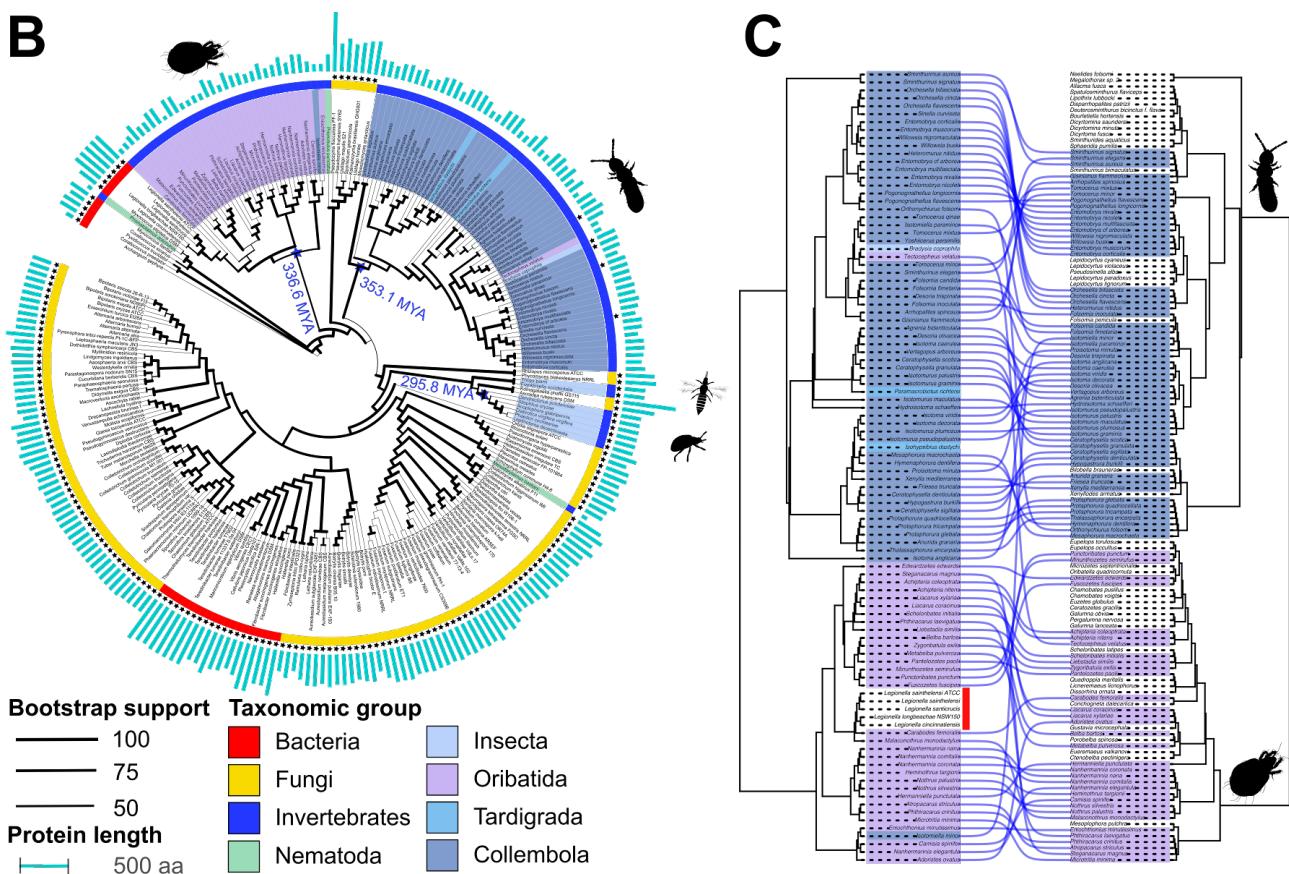
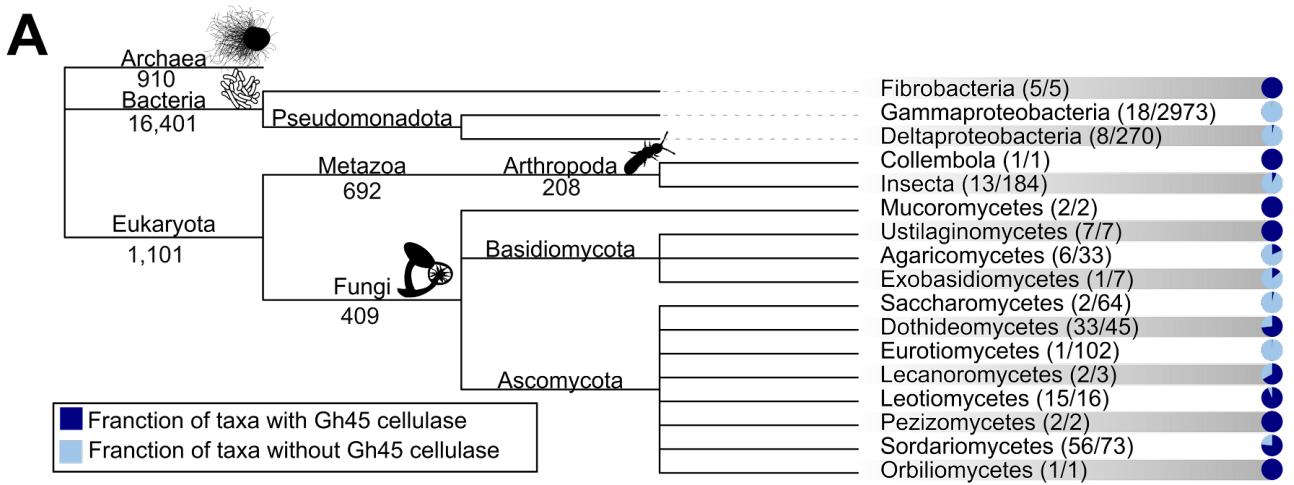
210 al., 2015), beetle (Busch et al., 2019) and Antarctic springtail (Hong et al., 2014) GH45
211 cellulases were all shown to be functional. Finally, orthologs with conserved domain
212 architectures were retained over hundreds of millions of years of evolution in springtails and
213 oribatids. This suggests little change has occurred in the trophic niche and position of
214 springtails and oribatid mites in soil food webs since their origin, further supported by the
215 presence of both taxonomic groups in the first fossil soils (Schaefer and Caruso, 2019;
216 Shear et al., 1984). Taken together, these are strong indications that GH45-type cellulases
217 in springtails and oribatids perform cellulose decomposition. Future work to experimentally
218 evaluate the functional properties of soil invertebrate cellulases (Song et al., 2017) should
219 consider all glycoside hydrolase genes, as gene duplication events may have led to
220 substrate diversification (Busch et al., 2019; Shin et al., 2022), with duplicates being able to
221 break down other polysaccharides like xyloglucan, mannans or xylan. While orthologs can
222 be identified bioinformatically, functional properties need to be confirmed experimentally by
223 expressing these enzymes heterologously, and test their substrate specificity on cellulose
224 and hemicellulose polysaccharides.

225 The ability for repeated HGT of GH45-type cellulases, and long term conservation in new
226 recipient taxa, may be related to several properties. First, cellulases, as secreted, gut-acting
227 enzymes (Fischer et al., 2013) do not depend on existing physiological pathways and their
228 regulation for proper functioning in the recipient organism. Second, such enzymes likely
229 reach their correct extracellular destination directly after the successful transfer of the
230 cellulase gene, its incorporation into the genome and its translation, because the signals of
231 protein export generally work independent of origin in most other taxa, even over vast
232 evolutionary distances (Clérigo et al., 2008). Third, the reaction catalyzed by cellulases
233 yields products that can serve as a beneficial fitness-relevant resource in any organism,
234 because the necessary downstream pathways are ubiquitously present. Supporting these
235 arguments, genes transferred horizontally are often secreted proteins (Savory et al., 2015;
236 Undheim and Jenner, 2021).

237 Horizontal gene transfer of cellulases, among other plant cell wall-degrading enzymes, is a
238 key process in the evolution of herbivory in arthropods (Wybouw et al., 2016). It resulted, for
239 example, in the massive radiation of Phytophaga, the most species-rich clade of beetles
240 (Busch et al., 2019), and in adaptation to lignocellulose-rich diets in crustaceans (King et al.,
241 2010). The long-term evolutionary preservation of GH45 genes suggests that cellulases
242 likely confer fitness benefits to soil invertebrates. These benefits may come from a direct use
243 of plant carbohydrate resources, although some theories imply that soil invertebrates are
244 limited rather by access to proteins, but not by access to carbohydrates. The ability to
245 degrade complex polysaccharides may also provide access to more nutritious, protein-rich
246 cytosols or microorganisms colonizing the inside of plant cells, such as saprotrophic fungi,
247 which are considered as major dietary components of both collembolans and oribatids
248 (Pollierer and Scheu, 2021). We expect that a taxonomically broad comparison of the
249 presence of cellulase genes with traits and trophic niches (Maraun et al., 2023; Potapov et
250 al., 2022) will provide insights into the functional ecology and evolution of soil invertebrates.

251 The wide-spread presence of cellulases in springtails and oribatid mites suggests that
252 invertebrates are independently capable of enzymatic lignocellulose decomposition, forming
253 a third evolutionarily and ecologically distinct group with such capability, in addition to
254 bacteria and fungi. This has important consequences for our understanding of soil food webs
255 and the soil carbon cycle. Fungi compared to bacteria are known to react differently to

256 environmental change such as experimental warming (Melillo et al., 2017) or habitat
257 degradation (Zhou et al., 2018). This results from key differences in life history strategies,
258 e.g. growth rates or nutrient use (Jansson and Hofmockel, 2020). Their differential reaction
259 to environmental change influences decomposition as distinct taxa determine the rate and
260 biochemical pathways of organic matter processing (Crowther et al., 2019). For example,
261 fungal-based, slow energy channels are more resistant to drought, limiting C and N losses
262 from agricultural soils (de Vries et al., 2012). Fungi accordingly contribute more to litter
263 decomposition than bacteria under drought conditions (Ullah et al., 2023). Soil invertebrates
264 react differently to environmental change compared to microorganisms (Sünnemann et al.,
265 2021). We hypothesize that global change has a more detrimental impact on decomposition
266 performed by soil invertebrates, given their lower effective population sizes and adaptive
267 elasticity (Lanfear et al., 2014; Pauls et al., 2013). It might be essential to consider these
268 differences for a better integration of below-ground processes into ecosystems models
269 (Chertov et al., 2017; Deckmyn et al., 2020; Filser et al., 2016) including global carbon
270 models (Friedlingstein et al., 2022), and for better predictions of soil carbon and nutrient
271 cycling.



273 **Fig. 1.** A) Abundance of GH45 cellulases in the three domains of life in NCBI RefSeq
274 genomes; B) Maximum likelihood phylogeny of the GH45-cellulase family. Branch lengths
275 are not drawn to scale, line weights indicate percent bootstrap support. Species represented
276 by a genome assembly in the NCBI RefSeq or GenBank databases are indicated by an
277 asterisk. Pictograms identify the four main soil invertebrate clades: oribatid mites, springtails,
278 thrips and beetles (clockwise). Bars mark GH45 protein length. Internal node labels provide
279 age estimates of the respective clades (Kumar et al., 2022); C) correspondence of oribatid
280 and springtail GH45 gene trees (left) and of a phylogenomic reconstruction of the species
281 phylogenies (right). D) Comparison of protein domain architecture of *F. candida* and *R.
282 solani* (left) and a 3D structure alignment of both proteins (right). The color indicates the
283 similarity of structures in the aligned proteins, with blue marking a high, and red marking low
284 correspondence. Lcr: low complexity region; Gh45: Pfam glyco-hydro 45 domain (PF02015).

285 Materials and Methods

286 Domain architectures of invertebrate GH45-type cellulases

287 Reviewed evidence exists for the presence of GH45-type cellulases in the antarctic springtail
288 (*Cryptopygus antarcticus*; Collembola, UniprotID D3GDK4) and the mustard beetle (*Phaedon
289 cochleariae*; Insecta, UniprotID O97401). Since these experimentally confirmed cellulases
290 harbor a Glyco-hydro 45 Pfam domain (PF02015) (Bankevich et al., 2012), we restricted our
291 analysis to cellulase orthologs that carry this Pfam domain. Pfam domains were annotated
292 with hmmscan from the HMMER package using Pfam v.32 using the default e-value cutoff of
293 0.01.

294 Genome Assembly pipeline

295 The genome assemblies provided by the MetalInvert Project (Bioproject ID: PRJNA758215)
296 cover a phylogenetically diverse set of soil-living invertebrates collected from the field or
297 obtained from cultures. Short read Illumina sequencing (300bp paired-end) with the
298 NovaSeq 6000 platform was done at Novogene Europe (Cambridge, UK), reads were
299 trimmed with Trimmomatic, human contaminating reads were filtered with Kraken2 (Wood
300 et al., 2019) and contig assembly was done with SPAdes (Bankevich et al., 2012). The
301 resulting contigs were taxonomically assigned with Blobtools2 (Challis et al., 2020) using the
302 NCBI non-redundant protein database as a reference, and only contigs with an assignment
303 to the phylum of the target species together with unassigned contigs were kept.
304 Redundancy reduction and scaffolding was done with Redundans (Pryszcz and Gabaldón,
305 2016), and genome assembly completeness was assessed with BUSCO (v 4.1.4) using the
306 precomputed metazoan (obd10) reference set. For our ortholog search, we selected
307 genomes with at least 50% BUSCO completeness (176 assemblies, Table S2).

308

309 RefSeq Gene Set collection

310 We downloaded gene sets for all 18,412 taxa represented in the NCBI RefSeq Genome
311 release 207 (O’Leary et al., 2016). The resulting taxon collections comprised 16,401 bacteria,
312 910 archaea, 409 fungi, 262 invertebrates and 430 vertebrates. The taxon list together with
313 the accession numbers are provided in Table S1.

314

315 Taxonomic assignment and contaminant detection
316 To rule out that fungal or bacterial contaminations of the underlying genome assemblies are
317 responsible for the animal cellulase orthologs (Steinegger and Salzberg, 2020), we
318 taxonomically classified each of the detected invertebrate orthologs. In brief, we used the
319 sequence as a query for a Diamond v.2.0.13 (Buchfink et al., 2015) search against the NCBI
320 non-redundant protein database (downloaded January 2022). From the resulting hit list, we
321 excluded the trivial hit against itself and then assigned the query sequence to the last
322 common ancestor of the taxa within a 10% bit score margin of the best hit (Huson et al.,
323 2016). A sequence was flagged as a putative contaminant if its taxonomic assignment was
324 not placed on the lineage from the species whose genome was analyzed to the root of the
325 tree of cellular life. This workflow is implemented into the software package taXaminer
326 (<https://github.com/BIONF/taxaminer>). This provided no evidence for a foreign origin of these
327 sequences (Table S3).

328 **Orthology-based phylogenetic profiles of fungal Gh45 cellulase**
329 Profile-based targeted ortholog searches in annotated gene sets were performed with fDOG
330 (Birikmen et al., 2021) using the GH45 cellulase of the fungus *Rhizoctonia solani* (NCBI
331 Accession XP_043186467.1) as the seed. For the training of the initial profile Hidden Markov
332 model we used the parameter `--minDist genus` and `--maxDist phylum` limiting the number of
333 training sequences to 6 (see Table S4 for more information). Candidate orthologs were
334 filtered for the presence of the Pfam glyco-hydro 45 domain (PF02015; see Table S5 for a
335 list of discarded orthologs). Ortholog search in the unannotated MetalInvert genome
336 assemblies were performed with the fDOG extension fDOG-Assembly. In brief, genomic
337 regions likely containing a GH45-type cellulase were identified with a tBLASTn search using
338 the consensus sequence included in the initial core gh45 core group from fDOG as query.
339 The hit region was extended by 500 nucleotides on either side and genes in the resulting
340 candidate genomic region were annotated with MetaEuk v5.34c21f2 (Levy Karin et al., 2020)
341 using the OMA database (release December 2021) (Nguyen et al., 2015) as the reference
342 database for the gene prediction. The corresponding protein sequences were then tested for
343 orthology using the routines of fDOG and afterwards features were annotated with FAS
344 (Dosch et al., 2023). The fDOG-assembly workflow is available from
345 https://github.com/BIONF/fDOG/tree/fdog_goes_assembly. The results from fDOG and
346 fDOG-Assembly were merged and visualized with PhyloProfile (Tran et al., 2018).

347 **Gh45 cellulase gene tree reconstruction**
348 To investigate the evolutionary history of the Gh45 cellulase, we used the identified
349 orthologs for a gene tree reconstruction. If the ortholog search obtained more than one co-
350 ortholog, we used the one that is most similar to the seed protein for the tree
351 reconstruction. Sequences were aligned with Muscle v3.8.1551 (Edgar, 2004) and alignment
352 columns comprising more than 50% gaps were removed with a custom perl script. The
353 resulting multiple sequence alignment was used as input for a maximum likelihood tree
354 reconstruction with IQ-TREE (Nguyen et al., 2015) v. 1.6.8. Branch support was assessed
355 with 1000 bootstrap replicates using the ultrafast bootstrap approach. The SH-aLRT branch
356 test was performed, and the optimal number of cores was automatically detected via IQ-
357 TREE (- nt AUTO). The gene tree was visualized with iTOL (Letunic and Bork, 2021). Animal
358 GH45-type cellulases are paraphyletic in this tree, and a topology test using the AU test
359 (Shimodaira and Hasegawa, 1999) confirmed that a tree with monophyletic animal
360 cellulases explained the data significantly worse (p-AU=9.2E-4). A second gene tree

361 containing all identified orthologs and co-orthologs was reconstructed with the same
362 approach. iTOL was used to prune and visualize the gene tree (containing all oribatids and
363 springtails) as well as to connect cellulase genes from the same species.

364

365 Phylogenies of springtails and oribatids

366 Phylogenies of springtail and oribatid organisms included in the MetalInvert project were
367 computed separately with a supermatrix approach. BUSCO version 5.4.2 (Simão et al.,
368 2015) with the precomputed BUSCO Arthropoda gene set (db10) was used to search for
369 orthologs in all genome assemblies. Multiple sequence alignments were computed with
370 MAFFT using local pairwise alignment and at maximum 1000 iterations (7.481)(Katoh and
371 Standley, 2013), trimmed with clipkit (1.3.0) (Steenwyk et al., 2020) and concatenated with
372 used FASconCAT-G (1.04) (Kück and Longo, 2014) into a supermatrix. Four phylogenetic
373 trees per taxon group were reconstructed with IQ-TREE (Nguyen et al., 2015). Branch
374 support was assessed with 1000 bootstrap replicates using the ultrafast bootstrap approach.
375 The best-fit model was Q.insect+F+R9 for oribatids, and Q.insect+F+R10 for springtails,
376 automatically chosen by ModelFinder according to BIC. The final consensus tree was
377 computed with splitstree (4.19.0)(Huson and Bryant, 2006) by summarizing the four IQ-
378 TREEs into a consensus tree. The consensus trees were outgroup-rooted using *Sarcoptes*
379 *scabiei* (GCA_020844145.1) for oribatids and *Machilis hrabei* (GCA_003456935.1),
380 *Drosophila albomicans* (GCA_009650485.1), and *Tyrophagus putrescentiae* from MetalInvert
381 for springtails as outgroups.

382 Correspondence of GH45 tree and phylogenies

383 The BUSCO-based phylogenies of springtails and oribatids were outgroup-rooted and
384 merged for the comparison with the Gh45 cellulase genetree. The tanglegram matches taxa
385 by links and was computed with R with packages phytools v1.4-0 (Revell, 2012) and castor
386 v1.7.6 (Louca and Doebeli, 2018).

387

388 3D structure comparison

389 The 3D structures of gh45 cellulase genes from *R.solani* and mustard beetle were retrieved
390 from precomputed predictions from UniProt (accession numbers A0A0B7FQX1 and
391 O97401). The structures were visualized and compared with VMD (Humphrey et al., 1996)
392 and the extensions MultiSeq (Roberts et al., 2006) in combination with the alignment tool
393 STAMP (Russell and Barton, 1992).

394

395 References

- 396 Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM, Nikolenko
397 SI, Pham S, Prjibelski AD, Pyshkin AV, Sirotnik AV, Vyahhi N, Tesler G, Alekseyev
398 MA, Pevzner PA. 2012. SPAdes: A New Genome Assembly Algorithm and Its
399 Applications to Single-Cell Sequencing. *J Comput Biol* **19**:455–477.
400 doi:10.1089/cmb.2012.0021
- 401 Berg M, Stoffer M, Heuvel H. 2004. Feeding guilds in Collembola based on digestive
402 enzymes. *Pedobiologia* **48**:589–601. doi:10.1016/j.pedobi.2004.07.006
- 403 Birikmen M, Bohnsack KE, Tran V, Somayaji S, Bohnsack MT, Ebersberger I. 2021. Tracing
404 Eukaryotic Ribosome Biogenesis Factors Into the Archaeal Domain Sheds Light on
405 the Evolution of Functional Complexity. *Front Microbiol* **12**:2598.
406 doi:10.3389/fmicb.2021.739000

- 407 Bradford MA, Veen GF (Ciska), Bonis A, Bradford EM, Classen AT, Cornelissen JHC,
408 Crowther TW, De Long JR, Freschet GT, Kardol P, Manrubia-Freixa M, Maynard DS,
409 Newman GS, Logtestijn RSP, Viketoft M, Wardle DA, Wieder WR, Wood SA, van der
410 Putten WH. 2017. A test of the hierarchical model of litter decomposition. *Nat Ecol
411 Evo* 1:1836–1845. doi:10.1038/s41559-017-0367-4
- 412 Briones MJI. 2018. The Serendipitous Value of Soil Fauna in Ecosystem Functioning: The
413 Unexplained Explained. *Front Environ Sci* 6.
- 414 Buchfink B, Xie C, Huson DH. 2015. Fast and sensitive protein alignment using DIAMOND.
415 *Nat Methods* 12:59–60. doi:10.1038/nmeth.3176
- 416 Busch A, Danchin EGJ, Pauchet Y. 2019. Functional diversification of horizontally acquired
417 glycoside hydrolase family 45 (GH45) proteins in Phytophaga beetles. *BMC Evol Biol*
418 19:100. doi:10.1186/s12862-019-1429-9
- 419 Challis R, Richards E, Rajan J, Cochrane G, Blaxter M. 2020. BlobToolKit – Interactive
420 Quality Assessment of Genome Assemblies. *G3 Genes Genomes Genet* 10:1361–
421 1374. doi:10.1534/g3.119.400908
- 422 Chang WH, Lai AG. 2018. Mixed evolutionary origins of endogenous biomass-
423 depolymerizing enzymes in animals. *BMC Genomics* 19:483. doi:10.1186/s12864-
424 018-4861-0
- 425 Chertov O, Komarov A, Shaw C, Bykhovets S, Frolov P, Shanin V, Grabarnik P, Pripitina I,
426 Zubkova E, Shashkov M. 2017. Romul_Hum—A model of soil organic matter
427 formation coupling with soil biota activity. II. Parameterisation of the soil food web
428 biota activity. *Ecol Model* 345:125–139. doi:10.1016/j.ecolmodel.2016.10.024
- 429 Clérigo EM, Maki JL, Giersach LM. 2008. Use of synthetic signal sequences to explore the
430 protein export machinery. *Pept Sci* 90:307–319. doi:10.1002/bip.20856
- 431 Cragg SM, Beckham GT, Bruce NC, Bugg TD, Distel DL, Dupree P, Etxabe AG, Goodell BS,
432 Jellison J, McGeehan JE, McQueen-Mason SJ, Schnorr K, Walton PH, Watts JE,
433 Zimmer M. 2015. Lignocellulose degradation mechanisms across the Tree of Life.
434 *Curr Opin Chem Biol, Energy • Mechanistic biology* 29:108–119.
435 doi:10.1016/j.cbpa.2015.10.018
- 436 Crowther TW, van den Hoogen J, Wan J, Mayes MA, Keiser AD, Mo L, Averill C, Maynard
437 DS. 2019. The global soil community and its influence on biogeochemistry. *Science*
438 365:eaav0550. doi:10.1126/science.aav0550
- 439 Davies GJ, Dodson GG, Hubbard RE, Tolley SP, Dauter Z, Wilson KS, Hjort C, Mikkelsen
440 JM, Rasmussen G, Schülein M. 1993. Structure and function of endoglucanase V.
441 *Nature* 365:362–364. doi:10.1038/365362a0
- 442 de Vries FT, Liiri ME, Bjørnlund L, Bowker MA, Christensen S, Setälä HM, Bardgett RD.
443 2012. Land use alters the resistance and resilience of soil food webs to drought. *Nat
444 Clim Change* 2:276–280. doi:10.1038/nclimate1368
- 445 Deckmyn G, Flores O, Mayer M, Domene X, Schnepf A, Kuka K, Looy KV, Rasse DP,
446 Briones MJI, Barot S, Berg M, Vanguelova E, Ostonen I, Vereecken H, Suz LM, Frey
447 B, Frossard A, Tiunov A, Frouz J, Grebenc T, Öpik M, Javaux M, Uvarov A,
448 Vindušková O, Krogh PH, Franklin O, Jiménez J, Yuste JC. 2020. KEYLINK: towards
449 a more integrative soil representation for inclusion in ecosystem scale models. I.
450 review and model concept. *PeerJ* 8:e9750. doi:10.7717/peerj.9750
- 451 Dosch J, Bergmann H, Tran V, Ebersberger I. 2023. FAS: assessing the similarity between
452 proteins using multi-layered feature architectures. *Bioinformatics* 39:btad226.
453 doi:10.1093/bioinformatics/btad226
- 454 Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high
455 throughput. *Nucleic Acids Res* 32:1792–1797. doi:10.1093/nar/gkh340
- 456 FAO, ITPS, GSBI, CBD, EC. 2020. State of knowledge of soil biodiversity - Status,
457 challenges and potentialities, Report 2020. Rome, Italy: FAO. doi:10.4060/cb1928en
- 458 Filser J, Faber JH, Tiunov AV, Brussaard L, Frouz J, De Deyn G, Uvarov AV, Berg MP,
459 Lavelle P, Loreau M, Wall DH, Querner P, Eijssackers H, Jiménez JJ. 2016. Soil
460 fauna: key to new carbon models. *SOIL* 2:565–582. doi:10.5194/soil-2-565-2016
- 461 Fischer BM, Meyer E, Maraun M. 2014. Positive correlation of trophic level and proportion of

- 462 sexual taxa of oribatid mites (Acari: Oribatida) in alpine soil systems. *Exp Appl Acarol*
463 **63**:465–479. doi:10.1007/s10493-014-9801-3
- 464 Fischer R, Ostafe R, Twyman RM. 2013. Cellulases from Insects In: Vilcinskas A, editor.
465 Yellow Biotechnology II: Insect Biotechnology in Plant Protection and Industry,
466 Advances in Biochemical Engineering/Biotechnology. Berlin, Heidelberg: Springer.
467 pp. 51–64. doi:10.1007/10_2013_206
- 468 Formenti G, Theissing K, Fernandes C, Bista I, Bombarely A, Bleidorn C, Ciofi C, Crottini
469 A, Godoy JA, Höglund J, Malukiewicz J, Mouton A, Oomen RA, Paez S, Palsbøll PJ,
470 Pampoulie C, Ruiz-López María J., Svardal H, Theofanopoulou C, de Vries J,
471 Waldvogel A-M, Zhang Guojie, Mazzoni CJ, Jarvis ED, Bálint M, Formenti G,
472 Theissing K, Fernandes C, Bista I, Bombarely A, Bleidorn C, Čiampor F, Ciofi C,
473 Crottini A, Godoy JA, Hoglund J, Malukiewicz J, Mouton A, Oomen RA, Paez S,
474 Palsbøll P, Pampoulie C, Ruiz-López María José, Svardal H, Theofanopoulou C, de
475 Vries J, Waldvogel A-M, Zhang Goujie, Mazzoni CJ, Jarvis E, Bálint M, Aghayan SA,
476 Alioto TS, Almudi I, Alvarez N, Alves PC, Amorim IR, Antunes A, Arribas P, Baldrian
477 P, Berg PR, Bertorelle G, Böhne A, Bonisoli-Alquati A, Boštjančić LL, Boussau B,
478 Breton CM, Buzan E, Campos PF, Carreras C, Castro LFi, Chueca LJ, Conti E,
479 Cook-Deegan R, Croll D, Cunha MV, Delsuc F, Dennis AB, Dimitrov D, Faria R,
480 Favre A, Fedrigo OD, Fernández R, Ficetola GF, Flot J-F, Gabaldón T, Galea Agius
481 DR, Gallo GR, Giani AM, Gilbert MTP, Grebenc T, Guschanski K, Guyot R, Hausdorf
482 B, Hawlitschek O, Heintzman PD, Heinze B, Hiller M, Husemann M, Iannucci A,
483 Irisarri I, Jakobsen KS, Jentoft S, Klinga P, Kloch A, Kratochwil CF, Kusche H,
484 Layton KKS, Leonard JA, Lerat E, Liti G, Manousaki T, Marques-Bonet T, Matos-
485 Maraví P, Matschiner M, Maumus F, Mc Cartney AM, Meiri S, Melo-Ferreira J,
486 Mengual X, Monaghan MT, Montagna M, Myslajek RW, Neiber MT, Nicolas V, Novo
487 M, Ozretić P, Palero F, Pârvulescu L, Pascual M, Paulo OS, Pavlek M, Pegueroles
488 C, Pellissier L, Pesole G, Primmer CR, Riesgo A, Rüber L, Rubolini D, Salvi D,
489 Seehausen O, Seidel M, Secomandi S, Studer B, Theodoridis S, Thines M, Urban L,
490 Vasemägi A, Vella A, Vella N, Vernes SC, Vernesi C, Vieites DR, Waterhouse RM,
491 Wheat CW, Wörheide G, Wurm Y, Zammit G. 2022. The era of reference genomes in
492 conservation genomics. *Trends Ecol Evol* **37**:197–202.
493 doi:10.1016/j.tree.2021.11.008
- 494 Friedlingstein P, O'Sullivan M, Jones MW, Andrew RM, Gregor L, Hauck J, Le Quéré C,
495 Luijkx IT, Olsen A, Peters GP, Peters W, Pongratz J, Schwingshackl C, Sitch S,
496 Canadell JG, Ciais P, Jackson RB, Alin SR, Alkama R, Arneth A, Arora VK, Bates
497 NR, Becker M, Bellouin N, Bittig HC, Bopp L, Chevallier F, Chini LP, Cronin M, Evans
498 W, Falk S, Feely RA, Gasser T, Gehlen M, Grätzalis T, Gloege L, Grassi G, Gruber
499 N, Gürses Ö, Harris I, Hefner M, Houghton RA, Hurt GC, Iida Y, Ilyina T, Jain AK,
500 Jersild A, Kadono K, Kato E, Kennedy D, Klein Goldewijk K, Knauer J, Korsbakken
501 JI, Landschützer P, Lefèvre N, Lindsay K, Liu J, Liu Z, Marland G, Mayot N, McGrath
502 MJ, Metzl N, Monacci NM, Munro DR, Nakaoaka S-I, Niwa Y, O'Brien K, Ono T,
503 Palmer PI, Pan N, Pierrot D, Pocock K, Poulter B, Resplandy L, Robertson E,
504 Rödenbeck C, Rodriguez C, Rosan TM, Schwingen J, Séférian R, Shutler JD,
505 Skjelvan I, Steinhoff T, Sun Q, Sutton AJ, Sweeney C, Takao S, Tanhua T, Tans PP,
506 Tian X, Tian H, Tilbrook B, Tsujino H, Tubiello F, van der Werf GR, Walker AP,
507 Wanninkhof R, Whitehead C, Willstrand Wranne A, Wright R, Yuan W, Yue C, Yue X,
508 Zaehle S, Zeng J, Zheng B. 2022. Global Carbon Budget 2022. *Earth Syst Sci Data*
509 **14**:4811–4900. doi:10.5194/essd-14-4811-2022
- 510 García-Palacios P, Maestre FT, Kattge J, Wall DH. 2013. Climate and litter quality differently
511 modulate the effects of soil fauna on litter decomposition across biomes. *Ecol Lett*
512 **16**:1045–1053. doi:10.1111/ele.12137
- 513 Greenslade P, Ireson JE. 1986. Collembola of the Southern Australian Culture Steppe and
514 Urban Environments: A Review of Their Pest Status and Key to Identification. *Aust J*
515 *Entomol* **25**:273–291. doi:10.1111/j.1440-6055.1986.tb01115.x
- 516 Griffiths HM, Ashton LA, Evans TA, Parr CL, Eggleton P. 2019. Termites can decompose

- 517 more than half of deadwood in tropical rainforest. *Curr Biol* **29**:R118–R119.
518 doi:10.1016/j.cub.2019.01.012
- 519 Griffiths HM, Ashton LA, Parr CL, Eggleton P. 2021. The impact of invertebrate
520 decomposers on plants and soil. *New Phytol* **231**:2142–2149. doi:10.1111/nph.17553
- 521 Han Z, Sieriebriennikov B, Susoy V, Lo W-S, Igreja C, Dong C, Berasategui A, Witte H,
522 Sommer RJ. 2022. Horizontally Acquired Cellulases Assist the Expansion of Dietary
523 Range in *Pristionchus* Nematodes. *Mol Biol Evol* **39**:msab370.
524 doi:10.1093/molbev/msab370
- 525 Hong SM, Sung HS, Kang MH, Kim C-G, Lee Y-H, Kim D-J, Lee JM, Kusakabe T. 2014.
526 Characterization of *Cryptopygus antarcticus* Endo- β -1,4-Glucanase from Bombyx
527 mori Expression Systems. *Mol Biotechnol* **56**:878–889. doi:10.1007/s12033-014-
528 9767-8
- 529 Humphrey W, Dalke A, Schulten K. 1996. VMD: visual molecular dynamics. *J Mol Graph*
530 **14**:33–38, 27–28. doi:10.1016/0263-7855(96)00018-5
- 531 Huson DH, Beier S, Flade I, Górska A, El-Hadidi M, Mitra S, Ruscheweyh H-J, Tappu R.
532 2016. MEGAN Community Edition - Interactive Exploration and Analysis of Large-
533 Scale Microbiome Sequencing Data. *PLOS Comput Biol* **12**:e1004957.
534 doi:10.1371/journal.pcbi.1004957
- 535 Huson DH, Bryant D. 2006. Application of Phylogenetic Networks in Evolutionary Studies.
536 *Mol Biol Evol* **23**:254–267. doi:10.1093/molbev/msj030
- 537 Jansson JK, Hofmockel KS. 2020. Soil microbiomes and climate change. *Nat Rev Microbiol*
538 **18**:35–46. doi:10.1038/s41579-019-0265-7
- 539 Joly F-X, Coq S, Coulis M, David J-F, Hättenschwiler S, Mueller CW, Prater I, Subke J-A.
540 2020. Detritivore conversion of litter into faeces accelerates organic matter turnover.
541 *Commun Biol* **3**:1–9. doi:10.1038/s42003-020-01392-4
- 542 Katoh K, Standley DM. 2013. MAFFT Multiple Sequence Alignment Software Version 7:
543 Improvements in Performance and Usability. *Mol Biol Evol* **30**:772–780.
544 doi:10.1093/molbev/mst010
- 545 Kern M, McGeehan JE, Streeter SD, Martin RNA, Besser K, Elias L, Eborall W, Malyon GP,
546 Payne CM, Himmel ME, Schnorr K, Beckham GT, Cragg SM, Bruce NC, McQueen-
547 Mason SJ. 2013. Structural characterization of a unique marine animal family 7
548 cellobiohydrolase suggests a mechanism of cellulase salt tolerance. *Proc Natl Acad
549 Sci* **110**:10189–10194. doi:10.1073/pnas.1301502110
- 550 King AJ, Cragg SM, Li Y, Dymond J, Guille MJ, Bowles DJ, Bruce NC, Graham IA,
551 McQueen-Mason SJ. 2010. Molecular insight into lignocellulose digestion by a
552 marine isopod in the absence of gut microbes. *Proc Natl Acad Sci* **107**:5345–5350.
553 doi:10.1073/pnas.0914228107
- 554 Kirsch R, Gramzow L, Theißen G, Siegfried BD, ffrench-Constant RH, Heckel DG, Pauchet
555 Y. 2014. Horizontal gene transfer and functional diversification of plant cell wall
556 degrading polygalacturonases: Key events in the evolution of herbivory in beetles.
557 *Insect Biochem Mol Biol* **52**:33–50. doi:10.1016/j.ibmb.2014.06.008
- 558 Kück P, Longo GC. 2014. FASconCAT-G: extensive functions for multiple sequence
559 alignment preparations concerning phylogenetic studies. *Front Zool* **11**:81.
560 doi:10.1186/s12983-014-0081-x
- 561 Kumar S, Suleski M, Craig JM, Kaspruwicz AE, Sanderford M, Li M, Stecher G, Hedges SB.
562 2022. TimeTree 5: An Expanded Resource for Species Divergence Times. *Mol Biol
563 Evol* **39**:msac174. doi:10.1093/molbev/msac174
- 564 Lanfear R, Kokko H, Eyre-Walker A. 2014. Population size and the rate of evolution. *Trends
565 Ecol Evol* **29**:33–41. doi:10.1016/j.tree.2013.09.009
- 566 Letunic I, Bork P. 2021. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic
567 tree display and annotation. *Nucleic Acids Res* **49**:W293–W296.
568 doi:10.1093/nar/gkab301
- 569 Levy Karin E, Mirdita M, Söding J. 2020. MetaEuk—sensitive, high-throughput gene
570 discovery, and annotation for large-scale eukaryotic metagenomics. *Microbiome*
571 **8**:48. doi:10.1186/s40168-020-00808-x

- 572 Lewin HA, Richards S, Lieberman Aiden E, Allende ML, Archibald JM, Bálint M, Barker KB,
573 Baumgartner B, Belov K, Bertorelle G, Blaxter ML, Cai J, Caperello ND, Carlson K,
574 Castilla-Rubio JC, Chaw S-M, Chen L, Childers AK, Coddington JA, Conde DA,
575 Corominas M, Crandall KA, Crawford AJ, DiPalma F, Durbin R, Ebenezer TE,
576 Edwards SV, Fedrigo O, Flicek P, Formenti G, Gibbs RA, Gilbert MTP, Goldstein
577 MM, Graves JM, Greely HT, Grigoriev IV, Hackett KJ, Hall N, Haussler D, Helgen
578 KM, Hogg CJ, Isobe S, Jakobsen KS, Janke A, Jarvis ED, Johnson WE, Jones SJM,
579 Karlsson EK, Kersey PJ, Kim J-H, Kress WJ, Kuraku S, Lawniczak MKN, Leebens-
580 Mack JH, Li X, Lindblad-Toh K, Liu X, Lopez JV, Marques-Bonet T, Mazard S, Mazet
581 JAK, Mazzoni CJ, Myers EW, O'Neill RJ, Paez S, Park H, Robinson GE, Roquet C,
582 Ryder OA, Sabir JSM, Shaffer HB, Shank TM, Sherkow JS, Soltis PS, Tang B,
583 Tedersoo L, Uliano-Silva M, Wang K, Wei X, Wetzer R, Wilson JL, Xu X, Yang H,
584 Yoder AD, Zhang G. 2022. The Earth BioGenome Project 2020: Starting the clock.
585 *Proc Natl Acad Sci U S A* **119**:e2115635118. doi:10.1073/pnas.2115635118
- 586 Louca S, Doebeli M. 2018. Efficient comparative phylogenetics on large trees. *Bioinformatics*
587 **34**:1053–1055. doi:10.1093/bioinformatics/btx701
- 588 Maraun M, Thomas T, Fast E, Treibert N, Caruso T, Schaefer I, Lu J-Z, Scheu S. 2023. New
589 perspectives on soil animal trophic ecology through the lens of C and N stable
590 isotope ratios of oribatid mites. *Soil Biol Biochem* **177**:108890.
591 doi:10.1016/j.soilbio.2022.108890
- 592 Melillo JM, Frey SD, DeAngelis KM, Werner WJ, Bernard MJ, Bowles FP, Pold G, Knorr MA,
593 Grandy AS. 2017. Long-term pattern and magnitude of soil carbon feedback to the
594 climate system in a warming world. *Science* **358**:101–105.
595 doi:10.1126/science.aan2874
- 596 Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: A Fast and Effective
597 Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol Biol Evol*
598 **32**:268–274. doi:10.1093/molbev/msu300
- 599 Nozaki M, Miura C, Tozawa Y, Miura T. 2009. The contribution of endogenous cellulase to
600 the cellulose digestion in the gut of earthworm (*Pheretima hilgendorfi*:
601 *Megascolecidae*). *Soil Biol Biochem* **41**:762–769. doi:10.1016/j.soilbio.2009.01.016
- 602 O'Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R, Rajput B, Robbertse B,
603 Smith-White B, Ako-Adjei D, Astashyn A, Badretdin A, Bao Y, Blinkova O, Brover V,
604 Chetvernin V, Choi J, Cox E, Ermolaeva O, Farrell CM, Goldfarb T, Gupta T, Haft D,
605 Hatcher E, Hlavina W, Joardar VS, Kodali VK, Li W, Maglott D, Masterson P,
606 McGarvey KM, Murphy MR, O'Neill K, Pujar S, Rangwala SH, Rausch D, Riddick LD,
607 Schoch C, Shkeda A, Storz SS, Sun H, Thibaud-Nissen F, Tolstoy I, Tully RE,
608 Vatsan AR, Wallin C, Webb D, Wu W, Landrum MJ, Kimchi A, Tatusova T, DiCuccio
609 M, Kitts P, Murphy TD, Pruitt KD. 2016. Reference sequence (RefSeq) database at
610 NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids
611 Res* **44**:D733–745. doi:10.1093/nar/gkv1189
- 612 Parey E, Louis A, Cabau C, Guiguen Y, Roest Croilius H, Berthelot C. 2020. Synteny-
613 Guided Resolution of Gene Trees Clarifies the Functional Impact of Whole-Genome
614 Duplications. *Mol Biol Evol* **37**:3324–3337. doi:10.1093/molbev/msaa149
- 615 Pauls SU, Nowak C, Bálint M, Pfenniger M. 2013. The impact of global climate change on
616 genetic diversity within populations and species. *Mol Ecol* **22**:925–946.
617 doi:<https://doi.org/10.1111/mec.12152>
- 618 Pausas JG, Bond WJ. 2020. On the Three Major Recycling Pathways in Terrestrial
619 Ecosystems. *Trends Ecol Evol* **35**:767–775. doi:10.1016/j.tree.2020.04.004
- 620 Phillips HRP, Guerra CA, Bartz MLC, Briones MJI, Brown G, Crowther TW, Ferlian O,
621 Gongalsky KB, Hoogen J van den, Krebs J, Orgiazzi A, Routh D, Schwarz B, Bach
622 EM, Bennett JM, Brose U, Decaëns T, König-Ries B, Loreau M, Mathieu J, Mulder C,
623 Putten WH van der, Ramirez KS, Rillig MC, Russell D, Rutgers M, Thakur MP, Vries
624 FT de, Wall DH, Wardle DA, Arai M, Ayuke FO, Baker GH, Beauséjour R, Bedano
625 JC, Birkhofer K, Blanchart E, Blossey B, Bolger T, Bradley RL, Callaham MA,
626 Capowicz Y, Caulfield ME, Choi A, Crotty FV, Crumsey JM, Dávalos A, Cosin DJD,

- 627 Dominguez A, Duhour AE, Eekeren N van, Emmerling C, Falco LB, Fernández R,
628 Fonte SJ, Fragoso C, Franco ALC, Fugère M, Fusilero AT, Gholami S, Gundale MJ,
629 López MG, Hackenberger DK, Hernández LM, Hishi T, Holdsworth AR, Holmstrup M,
630 Hopfensperger KN, Lwanga EH, Huhta V, Hurisso TT, Iannone BV, Iordache M,
631 Joschko M, Kaneko N, Kaniantska R, Keith AM, Kelly CA, Kernecker ML, Klaminder
632 J, Koné AW, Kooch Y, Kukkonen ST, Lalthanzara H, Lammel DR, Lebedev IM, Li Y,
633 Lidon JBJ, Lincoln NK, Loss SR, Marichal R, Matula R, Moos JH, Moreno G, Morón-
634 Ríos A, Muys B, Neirynck J, Norgrove L, Novo M, Nuutinen V, Nuzzo V, P95 MR,
635 Pansu J, Paudel S, Pérès G, Pérez-Camacho L, Piñeiro R, Ponge J-F, Rashid MI,
636 Rebollo S, Rodeiro-Iglesias J, Rodríguez MÁ, Roth AM, Rousseau GX, Rozen A,
637 Sayad E, Schaik L van, Scharenbroch BC, Schirrmann M, Schmidt O, Schröder B,
638 Seeber J, Shashkov MP, Singh J, Smith SM, Steinwandter M, Talavera JA, Trigo D,
639 Tsukamoto J, Valença AW de, Vanek SJ, Virto I, Wackett AA, Warren MW, Wehr
640 NH, Whalen JK, Wironen MB, Wolters V, Zenkova IV, Zhang W, Cameron EK,
641 Eisenhauer N. 2019. Global distribution of earthworm diversity. *Science* **366**:480–
642 485. doi:10.1126/science.aax4851
- 643 Pollierer MM, Scheu S. 2021. Stable isotopes of amino acids indicate that soil decomposer
644 microarthropods predominantly feed on saprotrophic fungi. *Ecosphere* **12**:e03425.
645 doi:10.1002/ecs2.3425
- 646 Post WM, Peng T-H, Emanuel WR, King AW, Dale VH, DeAngelis DL, others. 1990. The
647 global carbon cycle. *Am Sci* **78**:310–326.
- 648 Potapov A, Bellini BC, Chown SL, Deharveng L, Janssens F, Kováč L, Kuznetsova N, Ponge
649 J-F, Potapov M, Querner P. 2020. Towards a global synthesis of Collembola
650 knowledge: challenges and potential solutions. *Soil Org* **92**:161–188.
- 651 Potapov AM, Beaulieu F, Birkhofer K, Bluhm SL, Degtyarev MI, Devetter M, Goncharov AA,
652 Gongalsky KB, Klärner B, Korobushkin DI, Liebke DF, Maraun M, Mc Donnell RJ,
653 Pollierer MM, Schaefer I, Shrubbovich J, Semenyuk II, Sendra A, Tuma J, Tůmová M,
654 Vassilieva AB, Chen T-W, Geisen S, Schmidt O, Tiunov AV, Scheu S. 2022. Feeding
655 habits and multifunctional classification of soil-associated consumers from protists to
656 vertebrates. *Biol Rev* **97**:1057–1117. doi:10.1111/brv.12832
- 657 Potapov AM, Guerra CA, van den Hoogen J, Babenko A, Bellini BC, Berg MP, Chown SL,
658 Deharveng L, Kováč L, Kuznetsova NA, Ponge J-F, Potapov MB, Russell DJ,
659 Alexandre D, Alatalo JM, Arbea JI, Bandyopadhyaya I, Bernava V, Bokhorst S,
660 Bolger T, Castaño-Meneses G, Chauvat M, Chen T-W, Chomel M, Classen AT,
661 Cortet J, Čuchta P, Manuela de la Pedrosa A, Ferreira SSD, Fiera C, Filsier J,
662 Franken O, Fujii S, Koudji EG, Gao M, Gendreau-Berthiaume B, Gomez-Pamies DF,
663 Greve M, Tanya Handa I, Heiniger C, Holmstrup M, Homet P, Ivask M, Janion-
664 Scheepers C, Jochum M, Joimel S, Claudia S. Jorge B, Jucevica E, Ferlian O, Iuñes
665 de Oliveira Filho LC, Klauberg-Filho O, Baretta D, Krab EJ, Kuu A, de Lima ECA, Lin
666 D, Lindo Z, Liu A, Lu J-Z, Luciañez MJ, Marx MT, McCary MA, Minor MA, Nakamori
667 T, Negri I, Ochoa-Hueso R, Palacios-Vargas JG, Pollierer MM, Querner P,
668 Raschmanová N, Rashid MI, Raymond-Léonard LJ, Rousseau L, Saifutdinov RA,
669 Salmon S, Sayer EJ, Scheunemann N, Scholz C, Seeber J, Shveenkova YB,
670 Stebaeva SK, Sterzynska M, Sun X, Susanti WI, Taskaeva AA, Thakur MP, Tsiafouli
671 MA, Turnbull MS, Twala MN, Uvarov AV, Venier LA, Widenfalk LA, Winck BR,
672 Winkler D, Wu D, Xie Z, Yin R, Zeppelini D, Crowther TW, Eisenhauer N, Scheu S.
673 2023. Globally invariant metabolism but density-diversity mismatch in springtails. *Nat
674 Commun* **14**:674. doi:10.1038/s41467-023-36216-6
- 675 Pryszzcz LP, Gabaldón T. 2016. Redundans: an assembly pipeline for highly heterozygous
676 genomes. *Nucleic Acids Res* **44**:e113. doi:10.1093/nar/gkw294
- 677 Revell LJ. 2012. phytools: An R package for phylogenetic comparative biology (and other
678 things). *Methods Ecol Evol* **3**:217–223. doi:10.1111/j.2041-210X.2011.00169.x
- 679 Roberts E, Eargle J, Wright D, Luthey-Schulten Z. 2006. MultiSeq: unifying sequence and
680 structure data for evolutionary analysis. *BMC Bioinformatics* **7**:382.
681 doi:10.1186/1471-2105-7-382

- 682 Rosenberg Y, Bar-On YM, Fromm A, Ostikar M, Shoshany A, Giz O, Milo R. 2023. The
683 global biomass and number of terrestrial arthropods. *Sci Adv* **9**:eabq4049.
684 doi:10.1126/sciadv.abq4049
- 685 Russell RB, Barton GJ. 1992. Multiple protein sequence alignment from tertiary structure
686 comparison: assignment of global and residue confidence levels. *Proteins* **14**:309–
687 323. doi:10.1002/prot.340140216
- 688 Savory F, Leonard G, Richards TA. 2015. The Role of Horizontal Gene Transfer in the
689 Evolution of the Oomycetes. *PLOS Pathog* **11**:e1004805.
690 doi:10.1371/journal.ppat.1004805
- 691 Schaefer I, Caruso T. 2019. Oribatid mites show that soil food web complexity and close
692 aboveground-belowground linkages emerged in the early Paleozoic. *Commun Biol*
693 **2**:1–8. doi:10.1038/s42003-019-0628-7
- 694 Shear WA, Bonamo PM, Grierson JD, Rolfe WDI, Smith EL, Norton RA. 1984. Early Land
695 Animals in North America: Evidence from Devonian Age Arthropods from Gilboa,
696 New York. *Science* **224**:492–494. doi:10.1126/science.224.4648.492
- 697 Shelomi M, Heckel DG, Pauchet Y. 2016. Ancestral gene duplication enabled the evolution
698 of multifunctional cellulases in stick insects (Phasmatodea). *Insect Biochem Mol Biol*
699 **71**:1–11. doi:10.1016/j.ibmb.2016.02.003
- 700 Shelomi M, Watanabe H, Arakawa G. 2014. Endogenous cellulase enzymes in the stick
701 insect (Phasmatodea) gut. *J Insect Physiol* **60**:25–30.
702 doi:10.1016/j.jinsphys.2013.10.007
- 703 Shimodaira H, Hasegawa M. 1999. Multiple Comparisons of Log-Likelihoods with
704 Applications to Phylogenetic Inference. *Mol Biol Evol* **16**:1114.
705 doi:10.1093/oxfordjournals.molbev.a026201
- 706 Shin NR, Doucet D, Pauchet Y. 2022. Duplication of Horizontally Acquired GH5_2 Enzymes
707 Played a Central Role in the Evolution of Longhorned Beetles. *Mol Biol Evol*
708 **39**:msac128. doi:10.1093/molbev/msac128
- 709 Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015. BUSCO:
710 assessing genome assembly and annotation completeness with single-copy
711 orthologs. *Bioinformatics* **31**:3210–3212. doi:10.1093/bioinformatics/btv351
- 712 Song JM, Hong SK, An YJ, Kang MH, Hong KH, Lee Y-H, Cha S-S. 2017. Genetic and
713 Structural Characterization of a Thermo-Tolerant, Cold-Active, and Acidic Endo-β-
714 1,4-glucanase from Antarctic Springtail, Cryptopygus antarcticus. *J Agric Food Chem*
715 **65**:1630–1640. doi:10.1021/acs.jafc.6b05037
- 716 Steenwyk JL, Iii TJB, Li Y, Shen X-X, Rokas A. 2020. ClipKIT: A multiple sequence
717 alignment trimming software for accurate phylogenomic inference. *PLOS Biol*
718 **18**:e3001007. doi:10.1371/journal.pbio.3001007
- 719 Steinegger M, Salzberg SL. 2020. Terminating contamination: large-scale search identifies
720 more than 2,000,000 contaminated entries in GenBank. *Genome Biol* **21**:115.
721 doi:10.1186/s13059-020-02023-1
- 722 Sünnemann M, Siebert J, Reitz T, Schädler M, Yin R, Eisenhauer N. 2021. Combined
723 effects of land-use type and climate change on soil microbial activity and invertebrate
724 decomposer activity. *Agric Ecosyst Environ* **318**:107490.
725 doi:10.1016/j.agee.2021.107490
- 726 Tran N-V, Greshake Tzovaras B, Ebersberger I. 2018. PhyloProfile: dynamic visualization
727 and exploration of multi-layered phylogenetic profiles. *Bioinforma Oxf Engl* **34**:3041–
728 3043. doi:10.1093/bioinformatics/bty225
- 729 Ullah MR, Carrillo Y, Dijkstra FA. 2023. Relative contributions of fungi and bacteria to litter
730 decomposition under low and high soil moisture in an Australian grassland. *Appl Soil
731 Ecol* **182**:104737. doi:10.1016/j.apsoil.2022.104737
- 732 Undheim EAB, Jenner RA. 2021. Phylogenetic analyses suggest centipede venom arsenals
733 were repeatedly stocked by horizontal gene transfer. *Nat Commun* **12**:818.
734 doi:10.1038/s41467-021-21093-8
- 735 van den Hoogen J, Geisen S, Routh D, Ferris H, Traunspurger W, Wardle DA, de Goede
736 RGM, Adams BJ, Ahmad W, Andriuzzi WS, Bardgett RD, Bonkowski M, Campos-

- 737 Herrera R, Cares JE, Caruso T, de Brito Caixeta L, Chen X, Costa SR, Creamer R,
738 Mauro da Cunha Castro J, Dam M, Dijgal D, Escuer M, Griffiths BS, Gutiérrez C,
739 Hohberg K, Kalinkina D, Kardol P, Kergunteuil A, Korthals G, Krashevska V, Kudrin
740 AA, Li Q, Liang W, Magilton M, Marais M, Martín JAR, Matveeva E, Mayad EH,
741 Mulder C, Mullin P, Neilson R, Nguyen TAD, Nielsen UN, Okada H, Rius JEP, Pan K,
742 Peneva V, Pellissier L, Carlos Pereira da Silva J, Pitteloud C, Powers TO, Powers K,
743 Quist CW, Rasmann S, Moreno SS, Scheu S, Setälä H, Sushchuk A, Tiunov AV,
744 Trap J, van der Putten W, Vestergård M, Villenave C, Waeyenberge L, Wall DH,
745 Wilschut R, Wright DG, Yang J, Crowther TW. 2019. Soil nematode abundance and
746 functional group composition at a global scale. *Nature* **572**:194–198.
747 doi:10.1038/s41586-019-1418-6
- 748 Watanabe H, Noda H, Tokuda G, Lo N. 1998. A cellulase gene of termite origin. *Nature*
749 **394**:330–331. doi:10.1038/28527
- 750 Wood DE, Lu J, Langmead B. 2019. Improved metagenomic analysis with Kraken 2.
751 *Genome Biol* **20**:257. doi:10.1186/s13059-019-1891-0
- 752 Wybouw N, Pauchet Y, Heckel DG, Van Leeuwen T. 2016. Horizontal Gene Transfer
753 Contributes to the Evolution of Arthropod Herbivory. *Genome Biol Evol* **8**:1785–1801.
754 doi:10.1093/gbe/evw119
- 755 Zhou Z, Wang C, Luo Y. 2018. Effects of forest degradation on microbial communities and
756 soil carbon cycling: A global meta-analysis. *Glob Ecol Biogeogr* **27**:110–124.
757 doi:10.1111/geb.12663

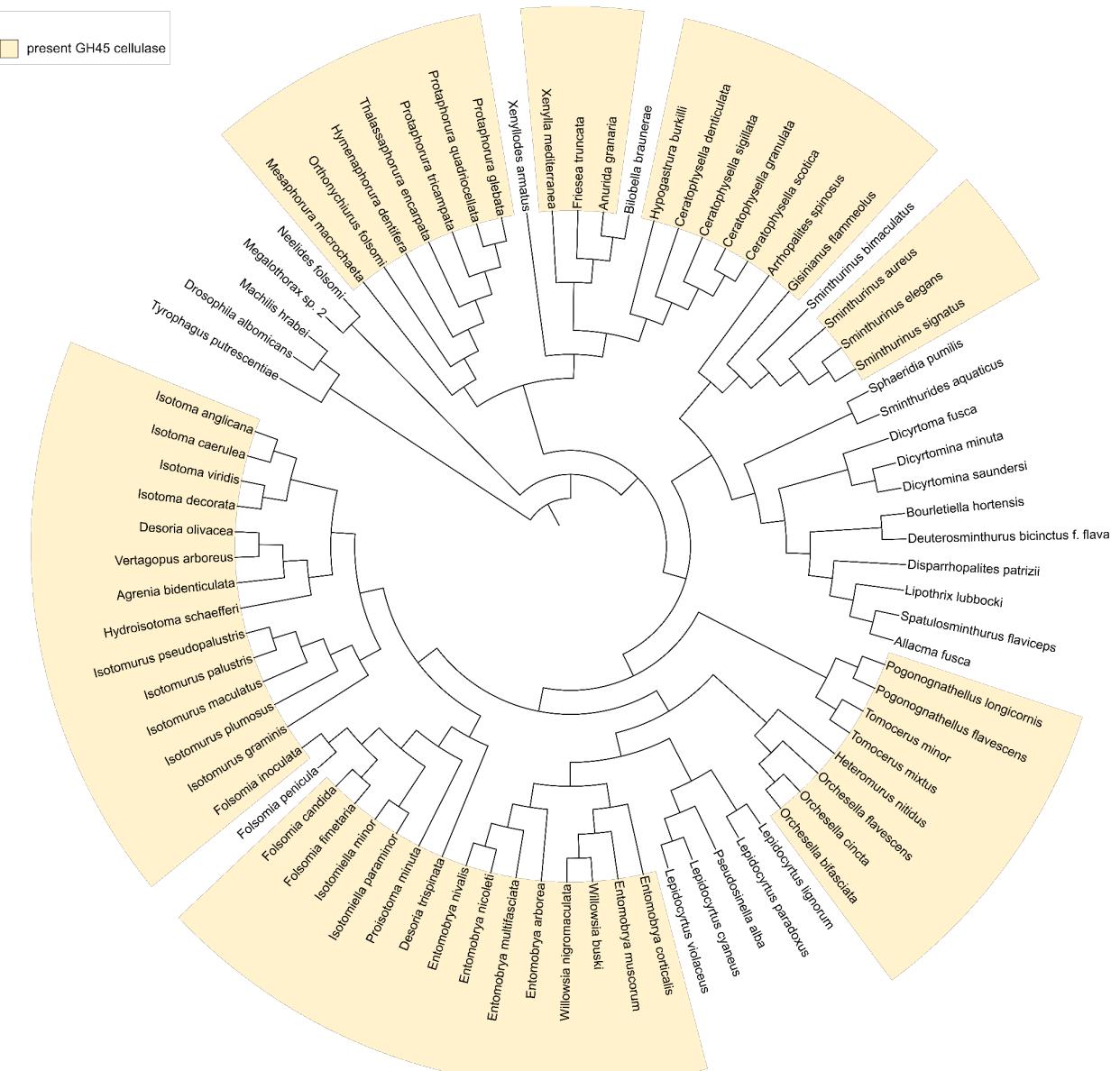
758 Acknowledgements

759 The research was funded by the Landes-Offensive zur Entwicklung Wissenschaftlich-
760 ökonomischer Exzellenz (LOEWE) Program of the Hessian Ministry of Higher Education,
761 Research, Science and the Arts through the LOEWE Centre for Translational Biodiversity
762 Genomics (LOEWE-TBG), and by the BA 4843/4-1 grant of the German Research
763 Foundation (DFG). MB, IE, HM, SS, IS, AP, MP, YP, GC conceptualized the study. IE, HM,
764 FA, MB, CS, JR, RL, KH, GC contributed methodology and materials. HM, FA, IE performed
765 the analyses and visualized results. IE supervised the analyses. MB and IE wrote the
766 original draft of the manuscript. HM, GC, TH, KH, RL, YP, MP, AP, JR, IS, SS, CS, IE and
767 MB reviewed and finalized the manuscript. All data needed to evaluate the conclusions in
768 the paper are present in the paper and/or the Supplementary Materials. Genome data will be
769 available via NCBI BioProject PRJNA758215. The authors declare that they have no
770 competing interests.

771 Supplementary Materials

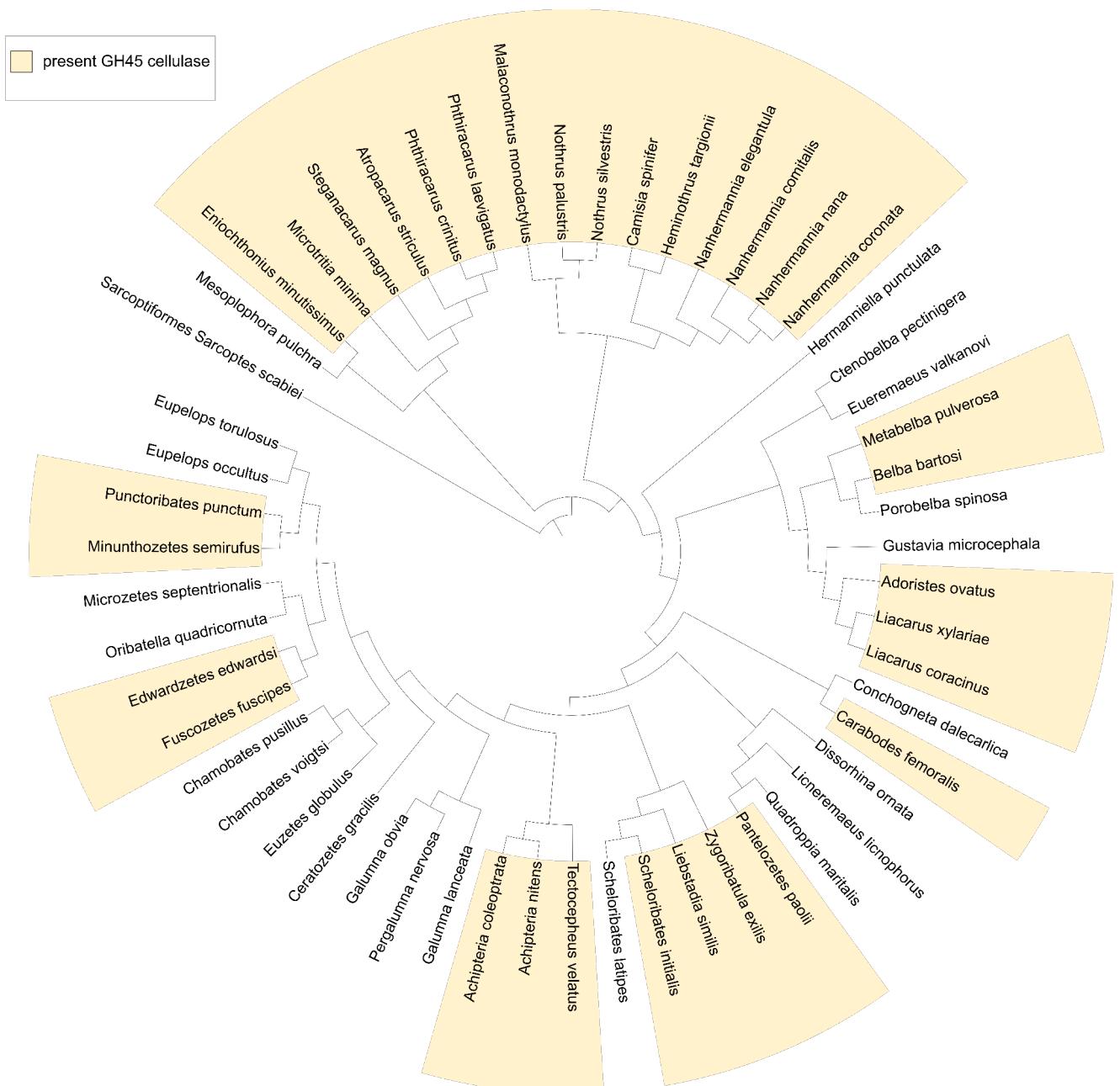
772 Supplementary Figures

773



774

Fig. S1. Maximum likelihood phylogeny of springtails. Species with GH45-cellulase are highlighted in yellow.



777

778 **Fig. S2.** Maximum likelihood phylogeny of oribatid mites. Species with GH45-cellulase are
779 highlighted in yellow.

780 **Fig. S3.** Maximum likelihood phylogeny of all GH45-cellulase genes identified in soil
781 invertebrates. Lines connect co-orthologs found in the same species. (uploaded separately)

782 Supplementary Tables

783
784 **Table S1. Screened RefSeq genomes.** Bacterial, archaeal, fungal, invertebrate and
785 vertebrate genomes screened for GH45 ortholog presence. (large table, uploaded
786 separately)

787
788 **Table S2. Screened MetalInvert genomes.** Soil invertebrate genomes of soil invertebrates
789 screened with fDOG-Assembly for GH45 ortholog presence.

libid	countgroup	Scientific name	taxid	busco completeness	Gh45 cellulase
a_17	Chilopoda	<i>Cryptops parisi</i>	173049	0,77	Absent
a_84	Chilopoda	<i>Geophilus carpophagus</i>	173285	0,82	Absent
P3_3	Chilopoda	<i>Geophilus flavus</i>	856749	0,72	Absent
a96	Chilopoda	<i>Geophilus truncorum</i>	173284	0,84	Absent
a_4	Chilopoda	<i>Haplophilus subterraneus</i>	173289	0,65	Absent
a_5	Chilopoda	<i>Henia vesuviana</i>	126936	0,8	Absent
a_21	Chilopoda	<i>Pachymerium ferrugineum</i>	115410	0,84	Absent
a_83	Chilopoda	<i>Stenotaenia linearis</i>	1569481	0,76	Absent
a42	Chilopoda	<i>Strigamia acuminata</i>	1255758	0,76	Absent
a39	Chilopoda	<i>Strigamia crassipes</i>	1428135	0,79	Absent
a_23	Chilopoda	<i>Strigamia transsilvanica</i>	1579475	0,7	Absent
a_66	Collembola	<i>Agrenia bidenticulata</i>	1933610	0,82	Present
MI_479	Collembola	<i>Allacma fusca</i>	39272	0,84	Absent
a_35	Collembola	<i>Anurida granaria</i>	187597	0,79	Present
a107	Collembola	<i>Arrhopalites spinosus</i>	187608	0,6	Present
a_42	Collembola	<i>Bilobella braunerae</i>	106916	0,57	Absent
a102	Collembola	<i>Bourletiella hortensis</i>	574228	0,76	Absent
P1_24	Collembola	<i>Ceratophysella denticulata</i>	928250	0,75	Present
a_34	Collembola	<i>Ceratophysella granulata</i>	1218962	0,78	Present
a52	Collembola	<i>Ceratophysella scotica</i>	187617	0,62	Present
a49	Collembola	<i>Ceratophysella sigillata</i>	1218965	0,66	Present
a72	Collembola	<i>Desoria olivacea</i>	370026	0,55	Present
MI_424	Collembola	<i>Desoria trispinata</i>	1184801	0,61	Present
		<i>Deuterosminthurus bicinctus f. flava</i>			
a3_16	Collembola	<i>Dicyrtoma fusca</i>	2041938	0,67	Absent
a_31	Collembola	<i>Dicyrtomina minuta</i>	1385863	0,8	Absent
a_32	Collembola	<i>Dicyrtomina saundersi</i>	1387116	0,77	Absent
P4_21	Collembola	<i>Disparrhopalites patrizii</i>	999999006	0,83	Absent
a_64	Collembola	<i>Entomobrya cf arborea</i>	30001	0,77	Present
a64	Collembola	<i>Entomobrya corticalis</i>	1503966	0,77	Present
a_33	Collembola	<i>Entomobrya muscorum</i>	2041940	0,53	Present
a_48	Collembola	<i>Entomobrya nicoleti</i>	2041941	0,7	Present
a71	Collembola	<i>Entomobrya nivalis</i>	1387109	0,66	Present
a108	Collembola	<i>Entomobrya Typ multifasciata</i>	247613	0,85	Present
a103	Collembola	<i>Folsomia candida</i>	158441	0,91	Present
a3_12	Collembola	<i>Folsomia fimetaria</i>	1387114	0,72	Present
a_38	Collembola	<i>Folsomia inoculata</i>	2041942	0,77	Present
MI_48	Collembola	<i>Folsomia penicula</i>	266765	0,59	Absent
a60	Collembola	<i>Friesea truncata</i>	187628	0,7	Present

a_40	Collembola	<i>Gisinianus flammeolus</i>	2449080	0,85	Present
P2_15	Collembola	<i>Heteromurus nitidus</i>	254095	0,81	Present
a98	Collembola	<i>Hydroisotoma schaefferi</i>	301519	0,85	Present
a58	Collembola	<i>Hymenaphorura dentifera</i>	999999008	0,72	Present
a_70	Collembola	<i>Hypogastrura burkilli</i>	1725397	0,71	Present
a105	Collembola	<i>Isotoma anglicana</i>	247611	0,66	Present
a_26	Collembola	<i>Isotoma caerulea</i>	308473	0,72	Present
a_71	Collembola	<i>Isotoma decorata</i>	57735	0,57	Present
P3_1	Collembola	<i>Isotoma viridis</i>	187635	0,75	Present
P1_25	Collembola	<i>Isotomiella minor</i>	370032	0,83	Present
a57	Collembola	<i>Isotomiella paraminor</i>	370031	0,58	Present
a_36	Collembola	<i>Isotomurus graminis</i>	1184803	0,72	Present
a_69	Collembola	<i>Isotomurus maculatus</i>	36143	0,66	Present
a97	Collembola	<i>Isotomurus palustris</i>	36144	0,81	Present
P1_26	Collembola	<i>Isotomurus plumosus</i>	1410395	0,53	Present
a_25	Collembola	<i>Isotomurus pseudopalustris</i>	36142	0,79	Present
a_53	Collembola	<i>Lepidocyrtus cyaneus</i>	247612	0,59	Absent
MI_426	Collembola	<i>Lepidocyrtus lignorum</i>	707889	0,59	Absent
a_27	Collembola	<i>Lepidocyrtus paradoxus</i>	49179	0,67	Absent
P2_2	Collembola	<i>Lepidocyrtus violaceus</i>	707891	0,81	Absent
a_46	Collembola	<i>Lipothrix lubbocki</i>	1387126	0,82	Absent
a4_24	Collembola	<i>Megalothorax sp. 2</i>	2340290	0,83	Absent
a68	Collembola	<i>Mesaphorura macrochaeta</i>	2651973	0,68	Present
MI_445	Collembola	<i>Neelides folsomi</i>	332381	0,75	Absent
a_29	Collembola	<i>Orchesella bifasciata</i>	576794	0,71	Present
a53	Collembola	<i>Orchesella cincta</i>	48709	0,53	Present
a_28	Collembola	<i>Orchesella flavescens</i>	48711	0,76	Present
a106	Collembola	<i>Orthonychiurus folsomi</i>	2581074	0,79	Present
a_49	Collembola	<i>Pogonognathellus flavescens</i>	511703	0,78	Present
a54	Collembola	<i>Pogonognathellus longicornis</i>	707266	0,84	Present
P2_17	Collembola	<i>Proisotoma minuta</i>	301521	0,84	Present
a61	Collembola	<i>Protaphorura glebata</i>	187683	0,79	Present
a_61	Collembola	<i>Protaphorura quadriocellata</i>	187683	0,72	Present
P2_19	Collembola	<i>Protaphorura tricampata</i>	187683	0,74	Present
a_41	Collembola	<i>Pseudosinella alba</i>	1302326	0,57	Absent
a104	Collembola	<i>Sminthurides aquaticus</i>	281415	0,8	Absent
a_39	Collembola	<i>Sminthurinus aureus</i>	1496267	0,84	Present
P4_22	Collembola	<i>Sminthurinus bimaculatus</i>	187699	0,81	Absent
P4_20	Collembola	<i>Sminthurinus elegans</i>	1190784	0,66	Present
a100	Collembola	<i>Sminthurinus signatus</i>	2584529	0,86	Present
a101	Collembola	<i>Spatulosminthurus flaviceps</i>	999999007	0,81	Absent
a3_8	Collembola	<i>Sphaeridia pumilis</i>	212016	0,87	Absent

P2_22	Collembola	<i>Thalassaphorura encarpata</i>	2583954	0,81 Present
P2_23	Collembola	<i>Tomocerus minor</i>	187706	0,58 Present
a55	Collembola	<i>Tomocerus mixtus</i>	58788	0,84 Present
a_43	Collembola	<i>Vertagopus arboreus</i>	2041954	0,62 Present
a_54	Collembola	<i>Willowsia buski</i>	1458441	0,77 Present
a_55	Collembola	<i>Willowsia nigromaculata</i>	1302335	0,8 Present
a99	Collembola	<i>Xenylla mediterranea</i>	2567731	0,78 Present
a62	Collembola	<i>Xenyllodes armatus</i>	187716	0,71 Absent
a94	Diplopoda	<i>Chordeuma sylvestre</i>	1569510	0,7 Absent
a_14	Diplopoda	<i>Cylindroiulus punctatus</i>	61981	0,68 Absent
a_96	Diplopoda	<i>Glomeris hexasticha</i>	1392624	0,66 Absent
a_10	Diplopoda	<i>Glomeris marginata</i>	62006	0,67 Absent
a_18	Diplopoda	<i>Julus scandinavius</i>	1008810	0,68 Absent
a_93	Diplopoda	<i>Julus scanicus</i>	541046	0,72 Absent
a_86	Diplopoda	<i>Kryphioiulus occultus</i>	1008825	0,72 Absent
a91	Diplopoda	<i>Megaphyllum sjaelandicum</i>	52423	0,7 Absent
a_95	Diplopoda	<i>Melogona broelemanni</i>	1147011	0,54 Absent
a95	Diplopoda	<i>Mycogona germanica</i>	999999013	0,67 Absent
a_13	Diplopoda	<i>Ommatoiulus sabulosus</i>	1008866	0,69 Absent
a_92	Diplopoda	<i>Ophyiulus pilosus</i>	118470	0,61 Absent
a_2	Diplopoda	<i>Polydesmus angustus</i>	1068628	0,72 Absent
a_15	Diplopoda	<i>Polydesmus complanatus</i>	510027	0,84 Absent
a_20	Diplopoda	<i>Proteroiulus fuscus</i>	88024	0,7 Absent
a_91	Diplopoda	<i>Rossiulus vilnensis</i>	999999014	0,72 Absent
a_12	Diplopoda	<i>Xestoiulus laeticollis</i>	1522044	0,74 Absent
MI_473	Enchytraeidae	<i>Cognettia cognettii</i>	1502715	0,58 Absent
a4_3	Enchytraeidae	<i>Enchytraeus crypticus</i>	913645	0,53 Absent
MI_474	Enchytraeidae	<i>Oconnorella tubifera</i>	913705	0,6 Absent
a35	Gamasina	<i>Phytoseiulus persimilis</i>	44414	0,8 Absent
a36	Gamasina	<i>Stratiolaelaps miles</i>	406085	0,83 Absent
a4_15	Nematoda	<i>Acrobeloides thornei</i>	96599	0,6 Absent
P2_7	Nematoda	<i>Aphelenchus avenae</i>	70226	0,59 Absent
MI_396	Nematoda	<i>Discolaimus major</i>	211252	0,61 Absent
a85	Nematoda	<i>Mesodorylaimus bastiani</i>	344383	0,55 Present
MI_395	Nematoda	<i>Panagrellus redivivus</i>	6233	0,55 Absent
a3_9	Nematoda	<i>Panagrolaimus detritophagus</i>	310956	0,55 Absent
P3_2	Nematoda	<i>Phasmarhabditis papillosa</i>	6243	0,64 Absent
a75	Nematoda	<i>Prionchulus punctatus</i>	293874	0,58 Present
D23	Nematoda	<i>Prismatolaimus dolichurus</i>	288633	0,53 Present
P2_4	Oribatida	<i>Achipteria coleoptrata</i>	229769	0,78 Present
P1_2	Oribatida	<i>Achipteria nitens</i>	229768	0,8 Present
a1	Oribatida	<i>Adoristes ovatus</i>	708363	0,72 Present

a2	Oribatida	<i>Atropacarus striculus</i>	229743	0,69 Present
a3	Oribatida	<i>Belba bartosi</i>	2241992	0,8 Present
a4	Oribatida	<i>Camisia spinifer</i>	198258	0,8 Present
MI_457	Oribatida	<i>Carabodes femoralis</i>	229793	0,74 Present
P1_32	Oribatida	<i>Ceratozetes gracilis</i>	1686620	0,77 Absent
a6	Oribatida	<i>Chamobates pusillus</i>	503572	0,75 Absent
a3_18	Oribatida	<i>Chamobates voigtsi</i>	198262	0,66 Absent
P4_1	Oribatida	<i>Conchogneta dalecarlica</i>	999999009	0,81 Absent
a7	Oribatida	<i>Ctenobelba pectinigera</i>	1401282	0,65 Absent
P1_12	Oribatida	<i>Dissorrhina ornata</i>	2202870	0,8 Absent
a10	Oribatida	<i>Edwardzetes edwardsi</i>	2202872	0,78 Present
D2	Oribatida	<i>Eniochthonius minutissimus</i>	229763	0,63 Present
a11	Oribatida	<i>Eueremaeus valkanovi</i>	1401269	0,74 Absent
P1_14	Oribatida	<i>Eupelops occultus</i>	2234141	0,77 Absent
a13	Oribatida	<i>Eupelops torulosus</i>	198282	0,77 Absent
P4_3	Oribatida	<i>Euzetes globulus</i>	334610	0,79 Absent
P1_4	Oribatida	<i>Fuscozetes fuscipes</i>	1686651	0,77 Present
P1_7	Oribatida	<i>Galumna lanceata</i>	229834	0,81 Absent
D3	Oribatida	<i>Galumna obvia</i>	885392	0,79 Absent
P1_8	Oribatida	<i>Gustavia microcephala</i>	1685391	0,65 Absent
MI_402	Oribatida	<i>Heminothrus targionii</i>	2664691	0,77 Present
		<i>Hermannilla punctulata</i> var. <i>septentrionalis</i>	885393	0,77 Absent
a15	Oribatida	<i>Liacarus coracinus</i>	198285	0,7 Present
a17	Oribatida	<i>Liacarus xyloiae</i>	198284	0,81 Present
MI_463	Oribatida	<i>Licneremaeus licnophorus</i>	999999011	0,77 Absent
P1_16	Oribatida	<i>Liebstadia similis</i>	1250587	0,8 Present
P4_5	Oribatida	<i>Malaconothrus monodactylus</i>	1797415	0,79 Present
MI_464	Oribatida	<i>Mesoplophora pulchra</i>	334620	0,79 Absent
MI_465	Oribatida	<i>Metabelba pulverosa</i>	229776	0,73 Present
a18	Oribatida	<i>Microtritia minima</i>	229747	0,75 Present
a19	Oribatida	<i>Microzetes septentrionalis</i>	999999012	0,74 Absent
P1_10	Oribatida	<i>Minunthozetes semirufus</i>	1979919	0,77 Present
P2_8	Oribatida	<i>Nanhermannia comitalis</i>	1979898	0,7 Present
D4	Oribatida	<i>Nanhermannia coronata</i> cf.	198290	0,75 Present
a20	Oribatida	<i>Nanhermannia elegantula</i>	66595	0,81 Present
MI_408	Oribatida	<i>Nanhermannia nana</i>	198291	0,78 Present
P1_11	Oribatida	<i>Nothrus palustris</i>	198293	0,74 Present
MI_467	Oribatida	<i>Nothrus silvestris</i>	66602	0,78 Present
P1_13	Oribatida	<i>Oribatella quadricornuta</i>	198298	0,75 Absent
P1_18	Oribatida	<i>Pantelozetes paolii</i>	1979943	0,74 Present
MI_412	Oribatida	<i>Pergalumna nervosa</i>	708370	0,78 Absent

a88	Oribatida	<i>Phthiracarus crinitus</i>	229740	0,82	Present
MI_492	Oribatida	<i>Phthiracarus laevigatus</i>	229740	0,81	Present
a22	Oribatida	<i>Porobelba spinosa</i>	2886740	0,79	Absent
a4_7	Oribatida	<i>Punc toribates punctum</i>	1720615	0,78	Present
P4_9	Oribatida	<i>Quadroppia maritalis</i>	1250640	0,7	Absent
a24	Oribatida	<i>Scheloribates initialis</i>	1979935	0,73	Present
D8	Oribatida	<i>Scheloribates latipes</i>	1979937	0,72	Absent
D7	Oribatida	<i>Steganacarus magnus</i>	52000	0,79	Present
a4_11	Oribatida	<i>Tectocepheus velatus</i>	229869	0,69	Present
P1_144	Oribatida	<i>Zygoribatula exilis</i>	1251916	0,76	Present
a4_20	Tardigrada	<i>Isohypsibius dastychi</i>	947160	0,64	Present
a3_3	Tardigrada	<i>Paramacrobiotus richtersi</i>	697321	0,57	Present

791

792

793 **Table S3. Metazoan genomes with GH45.** Taxonomic assignments of animal sequences
794 identified by fDOG as orthologs of cellulases with a GH45-type Pfam domain.
795

Scientific name	NCBI ID	Accession number	Protein ID	Gene ID	Taxonomic assignment
<i>Bradysia coprophila</i>	38358	GCF_014529535.1	XP_037026636.1	119067642	<i>Bradysia</i> <i>odoriphaga</i>
			XP_037050424.1	119084512	Protostomia
			XP_037027558.1	119068175	<i>Bradysia</i> <i>odoriphaga</i>
<i>Leptinotarsa decemlineata</i>	7539	GCF_000500325.1	XP_023016322.1	111505702	Chrysomelinae
			XP_023029513.1	111517551	<i>Gonioctena</i> <i>quinquepunctata</i>
			XP_023016323.1	111505703	Chrysomelinae
			XP_023029514.1	111517551	<i>Gonioctena</i> <i>quinquepunctata</i>
			XP_023022929.1	111511149	Chrysomelini
			XP_023016326.1	111505705	Chrysomelini
<i>Diabrotica virgifera</i>	50390	GCF_003013835.1	XP_028147313.1	114340743	Chrysomelidae
<i>virgifera</i>			XP_028139473.1	114333726	Chrysomelidae
			XP_028143849.1	114337572	Chrysomelidae
			XP_028147314.1	114340743	Chrysomelidae
<i>Anoplophora glabripennis</i>	217634	GCF_000390285.2	XP_018561275.1	108903540	<i>Anoplophora</i> <i>chinensis</i>

			XP_018561265.1	108903530	Lamiinae
<i>Dendroctonus ponderosae</i>	77166	GCF_000355655.1	XP_019754618.1	109533680	Dryophthorinae
			XP_019754620.1	109533682	Dryophthorinae
			XP_019771468.1	109545306	Cucujiformia
			XP_019754619.1	109533681	Dryophthorinae
			XP_019766961.1	109542255	Chrysomelinae
<i>Sitophilus oryzae</i>	7048	GCF_002938485.1	XP_030751361.1	115878892	Dryophthorinae
			XP_030747083.1	115875708	<i>Rhynchophorus ferrugineus</i>
<i>Thrips palmi</i>	161013	GCF_012932325.1	XP_034236588.1	117642458	<i>Frankliniella occidentalis</i>
<i>Frankliniella occidentalis</i>	133901	GCF_000697945.2	XP_026287984.1	113213214	<i>Thrips palmi</i>
			XP_026287985.1	113213214	<i>Thrips palmi</i>
			XP_026289264.1	113214189	<i>Thrips palmi</i>
<i>Folsomia candida</i>	158441	GCF_002217175.1	XP_021945337.1	110843646	Entomobryomorpha
			XP_021948187.1	110845935	Protostomia

796

797 **Table S4. Fungal GH45 inputs.** Core group of fungal species containing GH45 orthologs,
 798 computed by fDOG. This ortholog group was used as input for the final ortholog searches
 799 with fDOG and fDOG-Assembly by querying RefSeq genome assemblies for GH45
 800 presence.
 801

Scientific name	NCBI ID	Accession number RefSeq gene set	Gene IDs
<i>Rhizoctonia solani</i>	456999	GCF_016906535.1	XP_043186467.1
<i>Pleurotus ostreatus</i>	5322	GCF_014466165.1	XP_036634094.1
<i>Marasmius oreades</i>	181124	GCF_018924745.1	XP_043007043.1
<i>Pseudozyma flocculosa PF-1</i>	1277687	GCF_000417875.1	XP_007880076.1
<i>Kalmanozyma brasiliensis GHG001</i>	1365824	GCF_000497045.1	XP_016292070.1
<i>Ustilago maydis</i> 521	237631	GCF_000328475.2	XP_011388317.1

802

803 Table S5. **Excluded orthologs.** RefSeq candidate orthologs excluded from the final
 804 orthology inference results due to missing the Pfam domain characteristic of GH45
 805 cellulases.

Scientific name	NCBI ID	Class	RefSeq accession number	Gene ID from inferred ortholog
<i>Osmia lignaria</i>	473952	Insecta	GCF_012274295.1	XP_034173034.1
<i>Bombus terrestris</i>	30195	Insecta	GCF_000214255.1	XP_020720566.1 XP_012168823.1 XP_012168816.1 XP_012168808.1
<i>Colletes gigas</i>	935657	Insecta	GCF_013123115.1	XP_043266291.1 XP_043266289.1
<i>Dufourea novaeangliae</i>	178035	Insecta	GCF_001272555.1	XP_015435417.1
<i>Rhagoletis zephyri</i>	28612	Insecta	GCF_001687245.1	XP_017480646.1
<i>Streptomyces lacrimifluminis</i>	1500077	Actinomycetes	GCF_014646095.1	WP_189152206.1
<i>Actinoplanes globisporus DSM 43857</i>	1120949	Actinomycetes	GCF_000379645.1	WP_169516340.1
<i>Brevibacterium jeotgali</i>	1262550	Actinomycetes	GCF_007828155.1	WP_101587258.1
<i>Aneurinibacillus danicus</i>	267746	Bacilli	GCF_007991215.1	WP_146809708.1
<i>Paenibacillus thalictri</i>	2527873	Bacilli	GCF_004307995.1	WP_131011676.1
<i>Chondromyces apiculatus DSM 436</i>	1192034	Deltaproteobacteria	GCF_000601485.1	WP_197041519.1

806