1    Title

2    **Ecological selection for small microbial genomes along a temperate-to-thermal soil gradient**

3

4    Authors

5    Jackson W. Sorensen[1], Taylor K. Dunivin[1, 2], Tammy C. Tobin[3], and Ashley Shade[1,4]*

6      1.  Department of Microbiology and Molecular Genetics, Michigan State University, East

7          Lansing MI 48840 USA

8      2.  Environmental and Integrative Toxicological Sciences, Michigan State University, East

9          Lansing MI 48840

10      3.  Department of Biology, Susquehanna University, Selinsgrove, PA 17870 USA

11      4.  Department of Plant, Soil and Microbial Sciences; Program in Ecology, Evolutionary

12          Biology and Behavior; and the Plant Resilience Institute, Michigan State University, East

13          Lansing, MI 48840

14

15    *Materials and correspondence

16

17    Keywords

18    microbial ecology, genome size, genome reduction, thermophile, Centralia, coal seam fire,

19    metagenome, disturbance, extreme environment

20   Summary

21   Small bacterial and archaeal genomes provide insights into the minimal requirements for life[1]

22   and seem to be widespread on the microbial phylogenetic tree[2]. We know that evolutionary

23   processes, mainly selection and drift, can result in microbial genome reduction [3,4]. However, we

24   do not know the precise environmental pressures that constrain genome size in free-living

25   microorganisms. A study including isolates [5] has shown that bacteria with high optimum growth

26   temperatures, including thermophiles, often have small genomes [6]. It is unclear how well this

27   relationship may extend generally to microorganisms in nature [7,8], and in particular to those

28   microbes inhabiting complex and highly variable environments like soil [3,6,9]. To understand the

29   genomic traits of thermally-adapted microorganisms, here we investigated bacterial and

30   archaeal metagenomes from a 45°C gradient of temperate-to-thermal soils overlying the

31   ongoing Centralia, Pennsylvania (USA) coal seam fire. There was a strong relationship between

32   average genome size and temperature: hot soils had small genomes relative to ambient soils

33   (Pearson's r = -0.910, p < 0.001). There was also an inverse relationship between soil

34   temperature and cell size (Pearson's r = -0.65, p = 0.021), providing evidence that cell and

35   genome size in the wild are together constrained by temperature. Notably, hot soils had

36   different community structures than ambient soils, implicating ecological selection for thermo-

37   tolerant cells that had small genomes, rather than contemporary genome streamlining within

38   the local populations.  Hot soils notably lacked genes for described two-component regulatory

39   systems and antimicrobial production and resistance. Our work provides field evidence for the

40   inverse relationship between microbial genome size and temperature requirements in a

41   diverse, free-living community over a wide range of temperatures that support microbial life.

42    Our findings demonstrate that ecological selection for thermophiles and thermo-tolerant

43    microorganisms can result in smaller average genome sizes *in situ*, possibly because they have

44    small genomes reminiscent of a more ancestral state.

45
46

47    Main text

48        Genome streamlining is a reduction in genome size to increase cellular efficiency, and it

49    evolves by means of selection[3]. A comparative analysis of changes in microbial genomes sizes

50    with optimal growth temperature found a negative relationship that was independent of

51    phylogeny and environment [6]. This led to the conclusion that thermophiles are examples of

52    free-living microorganisms subject to genome streamlining [6,10,11]. These results were exciting

53    because they suggested that high temperature can select on genome size, providing insights

54    into environmental conditions that may propel efficiency. For the comparative analysis [6] and

55    cited studies therein, temperature optimum, genome size, 16S rRNA gene sequences, and

56    habitat were available for a curated collection 115 bacterial and archaeal isolates [5,12].  Given

57    biases of cultivation [13], an outstanding question was whether the relationship between growth

58    temperature and genome size would prove to be general for wild microbial communities.

59        Fortuitously, the fire-impact gradient at the Centralia ecosystem provides an

60    opportunity to investigate relationships between temperature and microbial genome traits.

61    Centralia, Pennsylvania is the site of a slow-burning, near-surface coal seam fire that ignited in

62    1962. The heat from the fire vents through overlying soils, causing surface soil temperatures to

63    reach as high as > 400°C [14], but more recently in the range of 40 - 75°C [15,16]. However, the soils

64    in Centralia were previously temperate, with no known exposure to prolonged high

3

65    temperatures. Therefore, Centralia offers an interesting model for the examining the eco-

66    evolution of microbial communities [17].

67         We recently used 16S rRNA gene amplicon sequencing to assess compositional changes

68    in Centralia soil microbial communities along an ambient-to-thermal temperature gradient [16].

69    Surface soils overlying the coal seam fire were collected to include soils that were hot from fire

70    ("fire-affected"), soils that were previously hot but had since recovered to ambient

71    temperatures ("recovered") and reference soils that had never been impacted by the fire. As

72    expected, fire-affected soils had starkly different community structure from ambient soils.

73    However, after the fire advanced, soils reasonably recovered towards reference community

74    structure. This suggested a considerable capacity of soil microbiomes for resilience, even after

75    exposure to a severe and unanticipated stressor, and prompted us to next ask what microbial

76    attributes underlay the observed changes in community structure in fire-affected soils.

77         Moving forward, we assessed average genome size along the Centralia fire gradient

78    (**Table S1**).  From twelve metagenomes (six fire-affected, five recovered, and one reference),

79    we used MicrobeCensus [18] to calculate average genome size across a soil temperature range of

80    45 °C. Average genome sizes were negatively and strongly correlated with temperature (**Figure**

81    **1A**, Pearson's $r = -0.910$, $p < 0.001$). In addition to MicrobeCensus, we used three other distinct

82    and complementary methods to assess changes in genome size with soil temperature and

83    found them all to be in agreement **(Figure S1).**  To the best of our knowledge, this is the first

84    report of decreases in genome size across an *in situ* temperature gradient that supports the

85    broad range of physiological requirements from mesophiles to thermophiles.

86       We next compared the average genome sizes estimated from Centralia metagenomes

87    to those from 22 publicly available soil metagenomes (**Figure 2, Table S2**).  Generally, hot soils

88    in Centralia had small genomes relative to other soils, while ambient soils in Centralia were

89    closer to the average size observed among this set.  Intriguingly, permafrost soils also harbored

90    small average genomes and were comparable to the hottest Centralia sites. These results

91    support comparably small genome sizes in Centralia soils and also provide a range of expected

92    soil genome sizes more generally.

93       It was hypothesized that small cells may be selected to attain minimal cellular

94    maintenance costs at high temperatures, and that small cells indirectly select for small

95    genomes [6]. Because we had microscope images from soil cell counts in Centralia[16], we re-

96    analyzed the images to extract size information.  We found that average cell sizes were also

97    negatively correlated with temperature (**Figure 1B**, Pearson's r = -0.65, p =0.021).  Accordingly,

98    cell size had a direct relationship with genome size (**Figure 1C**, Pearson's r = 0.64, p = 0.025).

99    These results agree with reported *in situ* relationships between cell size and temperature in

100    aquatic systems.  For example, an experiment investigating a 6°C increase in water temperature

101    confirmed that smaller cells with lower nucleotide content were selected at warmer

102    temperatures [7], providing support that even slight warming may enrich for microorganisms

103    with small genomes.  An observational study of marine microbial genome size along a

104    latitudinal gradient (10.7°C range) also supports this hypothesis [8].  Our results extend the cell

105    size-temperature trend to soils and also to a temperature range encompassing 45 °C.

106       To understand the selective outcomes of high temperature on the functions of these

107    small genomes, we next asked if there were functional genes that were characteristically

108    enriched or depleted with increasing temperature. We used shotgun metagenome annotations

109    from the KEGG module (KM) database [19]. KMs are groups of KEGG Orthologs (KOs) that

110    represent complexes, functional sets, metabolic pathways, or signatures. Eighty-one percent of

111    KOs detected in Centralia metagenomes were detected in all soils, and many patterns with

112    temperature were attributable to changes in normalized KO abundance rather than in KO

113    detection. In total, 284 (out of 541 detected; 52.50%) were correlated with temperature (**Figure**

114    **3**, **Table S3**).

115        Twenty-seven KMs were positively correlated with temperature (Pearson's R > 0.656,

116    false discovery rate adjusted p-value < 0.05;  **Figure 3A**). Specifically, dissimilatory sulfate

117    reduction (M00596), dissimilatory nitrate reduction (M00530) and denitrification (M00529)

118    were enriched in hot soils (**Figure 3A**, *cluster iii*; **Figure 4A**). These are anaerobic processes

119    aligned with known and expected environmental conditions in Centralia. Fire-affected soils

120    from active vents have higher moisture than reference and recovered soils (Pearson's r = 0.714,

121    p < 0.01), which likely promote inundated and anaerobic microhabitats therein. Prior work in

122    Centralia has indicated an importance of these metabolisms in hot soils, noting that sulfur,

123    sulfate, nitrate and ammonium were commonly elevated at vents [14,15]. These results also agree

124    with observations of thermophile metabolisms in other terrestrial and geothermal

125    environments, including a prevalence of denitrification and dissimilatory nitrate reduction [20,21],

126    highly active nitrogen cycles in hot springs [22], and increased dissimilatory organic sulfur

127    mineralization [23]. Notably, these anaerobic KMs grouped in their response patterns with several

128    archaeal proteins (**Figure 3A** *cluster iii*; Archaeal ribosome M00179, polymerase M00184, and

129    exosome M00390). We also observed an increase in Crenarchaeota in fire-affected soils [16], an

6

130    archaeal phylum that includes sulfate reducers [24].  Additional results describing patterns and

131    thresholds of KM enrichment with temperature are provided in Supporting Materials. Together,

132    these data suggest that the pathways enriched in small genomes from hot soils offer functions

133    attuned to the Centralia habitat.

134        Temperature was negatively correlated with 257 KMs (47.5% out of 541 total KMs

135    detected, Pearson's R < -0.6, false discovery rate adjusted p-value < 0.05; **Figure 3B**). In general,

136    these depleted KMs were detected across recovered soils and the reference soil**.** There were

137    two noteworthy categories of KMs that were consistently depleted in hot soils: antimicrobial

138    resistance and production and two component regulatory systems (**Figure 4B**).  Together, these

139    two KM categories comprised 32.7% of KMs negatively correlated with temperature (84 out of

140    257). This trend was striking, but we also note that some KMs belonging to these categories had

141    no relationships with temperature and that these KM categories were always detected in fire-

142    affected soils.

143        Thirty-nine modules for antimicrobial production and resistance mechanisms were

144    negatively correlated with temperature (**Figure 4B**)**,** which agrees with a prior analysis of

145    antibiotic resistance genes in this system [25]. Among these modules were resistance to

146    vancomycin, tetracycline, fluoroquinolone, aminoglycoside, nisin, erythromycin, streptomycin

147    and beta-lactam, and several multidrug efflux pumps. The small genomes of host-associated

148    symbionts often lack antimicrobial genes [26]. However, the *Pelagibacter* clade, which is a model

149    free-living population that has streamlined genomes, has a conserved multidrug transporter

150    across sequenced genomes [27]. It could be that thermophiles have fewer genes encoding

151    resistance to described antimicrobials, as evidenced by the challenges inherent in developing

7

152     specific selectable antibiotic resistance markers for thermophiles [28,29]. A related consideration

153     is that, like most databases, KEGG is biased towards genomes and annotations from fast-

154     growing mesophiles and may have missed annotation of under-described thermophile

155     antimicrobials. To clarify whether the observed decrease in antimicrobial production and

156     resistance was due to unannotated novelty or a true deficit of these functions in thermal sites,

157     annotation-independent methods could be used to identify antimicrobial-related biosynthetic

158     gene clusters from Centralia metagenomes [30,31]. In addition, functional screens of Centralia

159     isolates could be performed for antibiotic production and resistances. If there is a true deficit in

160     genes encoding antimicrobial production and resistance, it could be that the thermal conditions

161     present a strong environmental filter that reduces competition among the populations tolerant

162     of the heightened temperature. Our previous work reported decreased richness and

163     phylogenetic diversity fire-affected Centralia soils [16], suggesting that there is a smaller pool of

164     potential competitors inhabiting the hot soils.

165         Additionally, forty-nine detected two-component regulatory system modules were also

166     negatively correlated with temperature (Pearson's R < -0.6, **Figure 4B**). Two-component

167     systems consist of a sensor kinase and a response regulator and allow for transcriptional

168     responses to environmental stimuli [32]. This simple regulatory system allows bacteria to respond

169     to multiple stimuli: the involved genes duplicate, the sensors evolve sensitivity to additional

170     stimuli, and additional genes are transcribed [32,33]. Previous studies suggested that smaller

171     genomes have fewer regulatory components [34], and this relationship is often observed in

172     streamlined genomes [3,8]. Our results agree with observations of generally less regulation with

173     smaller genomes[4,11,27,35,36] and also suggest that thermophiles may have lower regulatory

174　　needs.  It has been proposed that thermophiles with "streamlined" genomes may be more

175　　likely to utilize global regulatory systems that mediate transcriptional responses to co-occurring

176　　environmental stimuli [11].  The degree of environmental variability is also predicted to influence

177　　the relative benefit an organism gains from investing in sensing its environment [37].  As a

178　　common case study in genome reduction, obligate endosymbionts are thought to have drifted

179　　towards small genomes in part because environmental conditions are stable and thus sensing

180　　requirements are minimal (e.g., [3] ).  Furthermore, in Centralia, seasonal temperature

181　　fluctuations in fire-affected soils are equivalent to those in ambient soils (**Figure S2**), providing

182　　evidence that the soils experience similar environmental stability in temperature, albeit at

183　　different ranges.  This suggests that small genomes are not necessarily conditional on very

184　　stable environments [3]. Future work should investigate whether two-component regulatory

185　　systems are consistently less prevalent among thermophiles, and, if so, whether their absence

186　　is reminiscent of an ancestral state.

187　　　　Our field study supports and reinforces cultivation-dependent observations that

188　　suggested bacteria and archaea with small genome sizes have higher growth temperatures [6].

189　　Because our study considers ecological section, as evidenced by the turnover in community

190　　membership between ambient and hot soils[16], these data indicate that environmental

191　　microorganisms with relatively higher temperature requirements also are likely to have small

192　　genomes and cell sizes. Surprisingly, it also suggests that microbial populations inhabiting

193　　complex environments, like soils, may generally reflect similar overarching traits in genome size

194　　as those observed in laboratory studies, which are necessarily biased towards fast-growing

195　　organisms that often are of medical, industrial, or agricultural interest (e.g., [38]). In addition, this

9

196    work expands upon previous reports of smaller genomes with higher temperatures [7,8] to

197    consider a range of *in situ* temperatures at which a variety of microbes compete in non-optimal

198    conditions. For example, we would expect mesophiles growing near their upper temperature

199    ranges and thermophiles growing near their lower temperature ranges to co-occur at some

200    sites in Centralia. Therefore, these results are relevant to the experiences of many wild

201    microorganisms that cope with dynamic environments.

202         Our results add evidence that supports both smaller genomes and cells, on average,

203    with higher temperatures but also offer a key point of distinction. Though the taxa enriched in

204    Centralia hot soils characteristically had smaller genomes and cells, there is no evidence for

205    contemporary genome streamlining in Centralia. Rather, we suspect that these thermo-tolerant

206    cells were resuscitated from the vast dormant pool in soil. This is supported by three lines of

207    evidence. First, there was turnover in community membership across hot and ambient

208    Centralia soils [16], providing evidence against contemporary streamlining within local lineages.

209    Second, there was striking comparability in average genome size of hot Centralia soils to

210    ancient permafrost soils, which largely contain an inactive and very old dormant pool.  Third,

211    many other studies have described thermophile persistence and resuscitation from non-

212    thermal environments, suggesting that these lineages are widespread but typically inactive [21,39–

213    43]. Therefore, we posit that Centralia small genomes are characteristic of an ancestral trait of

214    previously dormant thermophiles in the soil and not the outcome of genome streamlining.

215         In conclusion, we found a strong negative relationship between average microbial

216    genome size and temperature in Centralia soils along a mesophile-to-thermophile gradient,

217    spanning 45°C. We also found that cells were smaller in hot soils, supporting the hypothesis

10

218    that thermo-tolerant bacteria have smaller cell size, which indirectly selects for small genomes

219    [6]. By KEGG annotations, Centralia metagenomes at hot temperatures were best defined by

220    what they lacked rather than enriched modules of distinctive metabolisms. Specifically,

221    environmental sensing mechanisms, such as two-component regulatory systems, and

222    antimicrobial production and resistance mechanisms were in lower abundance in hot soils. In

223    addition, there were a few modules enriched at high temperatures that met expectations for

224    the hot anaerobic environment at active vents, including nitrogen and sulfur metabolism.  Our

225    results show that the relationship that was observed between growth temperature and

226    genome size for cultivable isolates also holds true in a complex, *in situ* microbial community

227    that inhabits a complex and variable soil environment. We suggest that, for thermo-tolerant

228    organisms, the relationship between temperature and genome size indicates the precursory

229    microbial condition of small genomes, reminiscent of ancient lineages, rather than

230    contemporary genome streamlining.

231

232    Materials and Methods

233    *DNA extraction and metagenome sequencing*

234    DNA for metagenome sequencing was manually extracted using a phenol chloroform extraction

235    [44] and then purified using the MoBio DNEasy PowerSoil Kit (MoBio, Solana Beach, CA, USA)

236    according the manufacturer's instructions. Total DNA sequencing was performed on all 12

237    samples by the Department of Energy's Joint Genome Institute (Community Science Project)

238    using an Illumina HiSeq 2500. Libraries were prepared with a targeted insert size of 270 base

239    pairs. Samples had between 19Gbp and 50Gbp of sequence data.  Additional methodology

240    details are provided in Supporting Materials.

241

242    *Quality control, assembly and annotation*

243    Assembly was performed by the Joint Genome Institute according to their standard operating

244    procedure (Supporting Materials). To use all sequencing data, we worked with assembled and

245    unassembled reads processed by Integrated Microbial Genomes (IMG) using their standard

246    annotation pipeline[45]. After comparing several annotation methods (Supporting Materials), we

247    chose to use the KEGG Orthology database for analyzing the Centralia data due to its inherent

248    structure and ability to integrate metabolic pathways. KEGG Ortholog (KO) abundances were

249    relativized to the median abundance in each site of a set of 36 single copy genes published

250    previously[46] (see Supporting Materials). One single copy gene (K01519) was an outlier in 7 out

251    of 12 samples as assessed by Grubb's test for outliers and removed. We analyzed patterns in

252    KEGG Modules (KMs)[19], a set of manually defined functional units made up of multiple KOs. KM

253    abundances were calculated based on the median abundance of their constituent KOs that

254    were present in the metagenomes. KMs were included in analysis if 50% or more of their

255    constituent KOs were identified in the dataset.  Approximately one third of the open reading

256    frames per sample were able to be annotated with KEGG (**Table S1**). As a caveat to the study,

257    unannotated open reading frames can result from erroneous reads and mis-assemblies but also

258    could be novel and or divergent genes critical for microbial processes. Thus, new annotations

259    could impact the overarching patterns described here.

260

261 *Average genome and cell size*

262 Average genome size was calculated from the quality filtered DNA sequences using

263 MicrobeCensus ("run_microbe_census.y –n 2000000"), which estimates average genome size

264 by calculating the percent of sampled reads that match to a set of single copy genes [18]. We also

265 used three additional methods to calculate average genome size (see Supporting Materials),

266 and all were in agreement in revealing the negative relationship between temperature and

267 average genome size. To calculate cell size, we re-analyzed microscope images previously used

268 to count microbial cells for community size quantifications in the same soils [16]. We hand-

269 curated a debris-free subset from the images and measured 44 - 910 cells from 3 - 9 replicate

270 fields for each soil. The major and minor axes of cells were measured using a FIJI macro in

271 ImageJ (Version: 2.0.0-rc-65/1.51s Build: 961c5f1b7f). We found that cell size range and

272 deviations (**Table S4**) were consistent with those previously reported [48].

273

274 *Comparisons with other soil metagenomes*

275 All metagenomic data sets for comparison were obtained from MG-RAST

276 ((http://metagenomics.anl.gov/). The MG-RAST database was searched with the following

277 criteria: material = soil, sequence type = shotgun, public = true. The resulting list of

278 metagenome data sets were ordered by number of base pairs (bp). Metagenomic data sets

279 with the most bp were included if they were sequenced using Illumina (to standardize

280 sequencing errors), had an available FASTQ file (for internal quality control), and contained <

281 30% low quality as determined by MG-RAST. Within high quality Illumina samples, priority for

282 inclusion was given to projects with multiple samples. When a project had multiple samples,

13

283     data sets with the greatest bp were selected. This search yielded 22 data sets from 12 locations

284     and five countries (**Table S2**). Sequences from MG-RAST data sets were quality checked using

285     FastQC (v0.11.3, [49] and quality controlled using the FASTX toolkit (fastq_quality_filter, "-Q33 -q

286     30 -p 50"). Average genome size for each dataset was calculated from the quality filtered DNA

287     sequences using MicrobeCensus with default parameters.

288

289     *Statistical analyses*

290          Statistics for the metagenome datasets were performed in the R environment for

291     statistical computing[50]. The stats package was used for calculating Pearson's correlations[50]. The

292     outliers package [51] was used for identifying outlying KOs. The ggplot2 package was used for

293     visualization[52]. Heat maps were created with heatmap2 from the gplots package[53].

294

295     *Data and workflows*

296     All analysis workflows are available on GitHub (ShadeLab/PAPER_SorensenInPrep).

297     Metagenome data are available on IMG under the GOLD Study ID GS0114513.

298

299     References

300     1.    Hutchison, C. A. *et al.* Design and synthesis of a minimal bacterial genome. *Science (80-. ).*
301           **351,** (2016).
302     2.    Hug, L. A. *et al.* A new view of the tree of life. *Nat. Microbiol.* **1,** 16048 (2016).
303     3.    Giovannoni, S. J., Thrash, J. C. & Temperton, B. Implications of streamlining theory for
304           microbial ecology. *ISME J.* **8,** 1553–1565 (2014).
305     4.    McCutcheon, J. P. & Moran, N. A. Extreme genome reduction in symbiotic bacteria.
306           *Nature Reviews Microbiology* **10,** 13–26 (2012).
307     5.    Vieira-Silva, S. & Rocha, E. P. C. The systemic imprint of growth and its uses in ecological
308           (meta)genomics. *PLoS Genet.* **6,** (2010).
309     6.    Sabath, N., Ferrada, E., Barve, A. & Wagner, A. Growth temperature and genome size in

310     bacteria are negatively correlated, suggesting genomic streamlining during thermal
311     adaptation. *Genome Biol. Evol.* **5,** 966–977 (2013).

312  7. Huete-Stauffer, T. M., Arandia-Gorostidi, N., Alonso-Sáez, L. & Morán, X. A. G.
313     Experimental warming decreases the average size and nucleic acid content of marine
314     bacterial communities. *Front. Microbiol.* **7,** (2016).

315  8. Swan, B. K. *et al.* Prevalent genome streamlining and latitudinal divergence of planktonic
316     bacteria in the surface ocean. *Proc. Natl. Acad. Sci.* **110,** 11463–11468 (2013).

317  9. Brewer, T. E., Handley, K. M., Carini, P., Gilbert, J. A. & Fierer, N. Genome reduction in an
318     abundant and ubiquitous soil bacterium 'Candidatus Udaeobacter copiosus'. *Nat.*
319     *Microbiol.* **2,** (2016).

320  10. Saha, D., Panda, A., Podder, S. & Ghosh, T. C. Overlapping genes: a new strategy of
321     thermophilic stress tolerance in prokaryotes. *Extremophiles* **19,** 345–353 (2015).

322  11. Wang, Q., Cen, Z. & Zhao, J. The Survival Mechanisms of Thermophiles at High
323     Temperatures: An Angle of Omics. *Physiology* **30,** 97–106 (2015).

324  12. Rocha, E. P. C. Codon usage bias from tRNA's point of view: Redundancy, specialization,
325     and efficient decoding for translation optimization. *Genome Res.* **14,** 2279–2286 (2004).

326  13. Rappé, M. S., Giovannoni, S. J., Rappe, M. S. & Giovannoni, S. J. The uncultured microbial
327     majority. *Annu. Rev. Microbiol.* **57,** 369–394 (2003).

328  14. Janzen, C. & Tobin-Janzen, T. Microbial communities in fire-affected soils. in *Microbiology*
329     *of Extreme Soils* 299–316 (Springer, 2008).

330  15. Tobin-Janzen, T. *et al.* Nitrogen changes and domain bacteria ribotype diversity in soils
331     overlying the Centralia, Pennsylvania underground coal mine fire. *Soil Sci.* **170,** (2005).

332  16. Lee, S.-H., Sorensen, J. W., Grady, K. L., Tobin, T. C. & Shade, A. Divergent extremes but
333     convergent recovery of bacterial and archaeal soil communities to an ongoing
334     subterranean coal mine fire. *ISME J.* **11,** 1447–1459 (2017).

335  17. Shade, A. Understanding microbiome stability in a changing world. *mSystems* **3,** e00157-
336     17 (2018).

337  18. Nayfach, S. & Pollard, K. S. Average genome size estimation improves comparative
338     metagenomics and sheds light on the functional ecology of the human microbiome.
339     *Genome Biol.* **16,** 51 (2015).

340  19. Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. & Morishima, K. KEGG: New
341     perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Res.* **45,** D353–
342     D361 (2017).

343  20. Torre, R. De, Dodsworth, J. A. & Hungate, B. Measuring Nitrification , Denitrification , and
344     Related Biomarkers in Terrestrial Geothermal Ecosystems. **486,** (2011).

345  21. Marchant, R. *et al.* Thermophilic bacteria in cool temperate soils: Are they metabolically
346     active or continually added by global atmospheric transport? *Appl. Microbiol. Biotechnol.*
347     **78,** 841–852 (2008).

348  22. Reigstad, L. J. *et al.* Nitrification in terrestrial hot springs of Iceland and Kamchatka. *FEMS*
349     *Microbiol. Ecol.* **64,** 167–174 (2008).

350  23. Santana, M., Gonzalez, J. & Clara, M. Inferring pathways leading to organic-sulfur
351     mineralization in the Bacillales. *Crit. Rev. Microbiol.* **42,** 31–45 (2016).

352  24. Itoh, T., Suzuki, K., Sanchez, P. C. & Nakase, T. Caldivirga maquilingensis gen. nov., sp.
353     nov., a new genus of rod-shaped crenarchaeote isolated from a hot spring in the

354   Philippines. *Int J Syst Bacteriol* **49,** 1157–1163 (1999).

355 25. Dunivin, T. K. & Shade, A. Community structure explains antibiotic resistance gene
356   dynamics over a temperature gradient in soil. *FEMS Microbiol. Ecol.* **fiy016,** (2018).

357 26. Gao, Z. M. *et al.* Symbiotic Adaptation Drives Genome Streamlining of the Cyanobacterial
358   Sponge Symbiont 'Candidatus Synechococcus spongiarum'. *MBio* **5,** 1–11 (2014).

359 27. Grote, J. *et al.* Streamlining and Core Genome Conservation among Highly Divergent
360   Members of the SAR11 Clade. *MBio* **3,** 1–13 (2012).

361 28. Brouns, S. J. J. *et al.* Engineering a Selectable Marker for Hyperthermophiles. *J. Biol.*
362   *Chem.* **280,** 11422–11431 (2005).

363 29. Hoseki, J., Yano, T., Koyama, Y. & Kuramitsu, S. Directed Evolution of Thermostable
364   Kanamycin-Resistance Gene : A Convenient Selection Marker for Thermus thermophilus
365   1. **956,** 951–956 (1999).

366 30. Weber, T. *et al.* antiSMASH 3 . 0 — a comprehensive resource for the genome mining of
367   biosynthetic gene clusters. **43,** 237–243 (2017).

368 31. Hadjithomas, M. *et al.* IMG-ABC : A Knowledge Base To Fuel Discovery of Biosynthetic
369   Gene Clusters and Novel Secondary Metabolites. *MBio* **6,** 1–10 (2015).

370 32. Hoch, J. A. Two-component and phosphorelay signal transduction. 165–170 (2000).

371 33. Whitworth, D. E. & Cock, Æ. P. J. A. Evolution of prokaryotic two-component systems :
372   insights from comparative genomics. 459–466 (2009). doi:10.1007/s00726-009-0259-2

373 34. Ranea, J. A. G., Grant, A., Thornton, J. M. & Orengo, C. A. Microeconomic principles
374   explain an optimal genome size in bacteria. **21,** 21–25 (2005).

375 35. Moran, N. A. Microbial minimalism: Genome reduction in bacterial pathogens. *Cell* **108,**
376   583–586 (2002).

377 36. Yus, E. *et al.* Impact of Genome Reduction on Bacterial Metabolism and Its Regulation.
378   *Science (80-. ).* **326,** 1263–1268 (2009).

379 37. Kussell, E. & Leibler, S. Phenotypic diversity, population growth, and information in
380   fluctuating environments. *Science* **309,** 2075–2078 (2005).

381 38. Gweon, H. S., Bailey, M. J. & Read, D. S. Assessment of the bimodality in the distribution
382   of bacterial genome sizes. *ISME J.* **11,** 821–824 (2017).

383 39. Rahman, T. J., Marchant, R. & Banat, I. M. Distribution and molecular investigation of
384   highly thermophilic bacteria associated with cool soil environments. *Biochem. Soc. Trans.*
385   **32,** 209–213 (2004).

386 40. Portillo, M. C., Santana, M. & Gonzalez, J. M. Presence and potential role of thermophilic
387   bacteria in temperate terrestrial environments. *Naturwissenschaften* **99,** 43–53 (2012).

388 41. Perfumo, A. & Marchant, R. Global transport of thermophilic bacteria in atmospheric
389   dust. *Environ. Microbiol. Rep.* **2,** 333–339 (2010).

390 42. Hubert, C. *et al.* A Constant Flux of Diverse Thermophilic Bacteria into the Cold Arctic
391   Seabed. *Science (80-. ).* **325,** 1541–1544 (2009).

392 43. Müller, A. L. *et al.* Endospores of thermophilic bacteria as tracers of microbial dispersal
393   by ocean currents. *ISME J.* **8,** 1153–65 (2014).

394 44. Cho, J.-C., Lee, D.-H., Cho, Y.-C., Cho, J.-C. & Kim, S.-J. Direct Extraction of DNA from Soil
395   for Amplification of 16S rRNA Gene Sequences by Polymerase Chain Reaction. *J.*
396   *Microbiology* 229–235 (2006).

397 45. Huntemann, M. *et al.* The standard operating procedure of the DOE-JGI Metagenome

398          Annotation Pipeline (MAP v.4). *Stand. Genomic Sci.* **11,** (2016).

399   46.   He, S. *et al.* Patterns in wetland microbial community composition and functional gene
400        repertoire associated with methane emissions. *MBio* **6,** e00066-15 (2015).

401   47.   Grubbs, F. E. Sample criteria for testing outlying observations. *Ann. Math. Stat.* 27–58
402        (1950).

403   48.   Balkwill, D. L. & Casida, L. E. Microflora of soil as viewed by freeze etching. *J. Bacteriol.*
404        **114,** 1319–1327 (1973).

405   49.   Andrews, S. FastQC: a quality control tool for high throughput sequence data. (2010).

406   50.   R Core Team. R: A Language and Environment for Statistical Computing. (2017).

407   51.   Komsta, L. outliers: Tests for outliers. *R Packag. version 0.14* . http://CRAN.R-
408        project.org/package=outliers (2011). doi:doi:10.1201/9780203910894.ch6

409   52.   Wickham, H. *ggplot2: Elegant graphics for data analysis*. (Springer-Verlag, 2009).

410   53.   Warnes, G. R. *et al.* gplots: Various R Programming Tools for Plotting Data. (2016).

411

412


413 Acknowledgements

421


422 Contributions

423 AS and TCT conceived the study and conducted field work.  JWS and TKD performed analyses,

424 with direction and oversight by AS.  JWS, AS and TKD contributed writing.  All authors discussed

425 results, and commented on and edited the manuscript.

426

427    Competing financial interests

428    The authors declare no competing financial interests.

429

430    Figure Legends

431    Figure 1. Changes in average genome and cell sizes across the soil temperature gradient in

432    Centralia. (A) Average genome size in each metagenome was calculated using MicrobeCensus

433    and plotted against site temperature.  (B) Average cell length was measured from 44-910 cells

434    from 3-9 replicate fields for each soil and plotted against soil temperature.  (C) Average genome

435    size had a direct relationship with average cell size.

436

437    Figure 2.  Average genome size in soil metagenomes, estimated using MicrobeCensus.[18]

438    Samples are ordered by average genome size and colored by sample location.

439

440    Figure 3. Heatmap of KEGG modules correlated with temperature (false discovery rate adjusted

441    p-value < 0.05).  Modules (rows) are centered and standardized across Centralia metagenomes

442    (columns), with warm colors showing relative enrichment and cool colors showing relative

443    depletion. Modules with significant relationships with temperature are shown. Sites are

444    arranged by increasing temperature from left to right. (A) 27 KEGG modules were positively

445    correlated with temperature (Pearson's R range = 0.646 to 0.933). (B) 257 KEGG modules were

446    negatively correlated with temperature (Pearson's R range = -0.642 to -0.925). A third of the

447    KEGG modules negatively correlated with temperature were either two-component regulatory

448     systems (TCRS, blue dendrogram tips), antimicrobial resistance or production (ARP, gray tips),

449     or both (black tips). Note differences in color gradient ranges across panels A and B.

450

451     Figure 4. KEGG modules that had notable enrichments or depletions with temperature. (A) The

452     median abundances of KEGG modules for denitrification (red), dissimilatory nitrate reduction

453     (green) and dissimilatory sulfate reduction (blue) were all positively correlated with

454     temperature. (B) Pearson's correlation values for all detected modules classified as antibiotic

455     resistance and production (gray density, n = 62 detected modules) or two-component

456     regulatory systems (blue density, n = 89 detected modules).  The black vertical line

457     distinguishes correlation values that are significant at a false discovery rate adjusted p-value <

458     0.05 (left), and all of these had a strong and negatively relationship with temperature. In total,

459     there were 39 antimicrobial resistance and production modules and 49 two-component

460     regulatory system modules that significantly decreased with temperature.

461

462     Supporting Figures

463     Figure S1.  Complementary methods used to assess changes in average genome size across the

464     soil temperature gradient in Centralia. (A) Odds ratios were calculated for 35 single-copy gene

465     KEGG Orthologs in each site and plotted against site temperature. Reported correlation is

466     between all single copy gene odds ratios and temperature, and all p < 0.001. (B) Average

467     genome size in each site was calculated based on phylum level abundances from 16S rRNA gene

468     amplicon data, using weighted average genome sizes of each phylum present in JGI IMG

469    (accessed 19 June 2017, correlation p < 0.001). (C) Average MAG size at each site was calculated

470    based on presence/absence of 104 MAGs (correlation p = 0.029).

471

472    Figure S2.  Annual temperature fluctuations at three fire-affected (circles) and two ambient

473    (triangles) Centralia sites, measured using *in situ* temperature loggers (HOBOs) that were buried

474    5 - 10 cm below the surface. Temperature loggers were deployed after the soils were collected

475    for this study.

476

477    Supporting Tables

478    Table S1. Sequence summary information for Centralia metagenomes.  Soils were collected 03-
479    07 October 2014. Asterisks indicate that the site was actively venting at the time of soil
480    collection.
481    Table S2.  MG-RAST metadata for soil metagenomes used in this study.
482    Table S3.  KEGG Modules significantly correlated with temperature (false-discovery-rate
483    adjusted p-value <0.05)
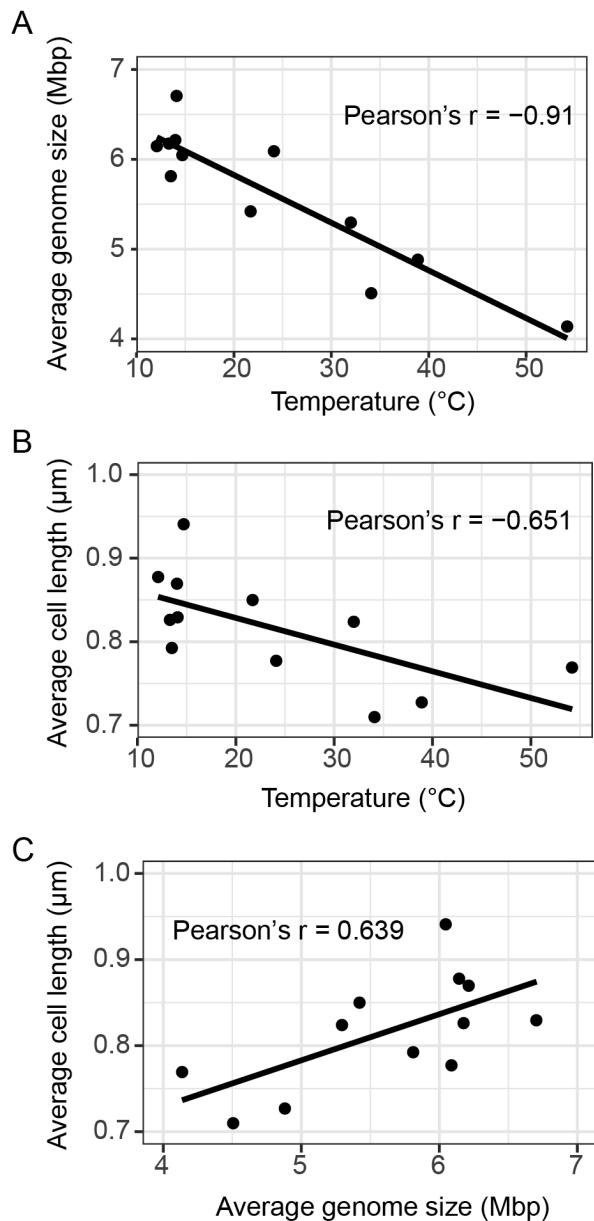484    Table S4. Cell size measurements from microscope images.
485    Table S5. Single-copy KEGG Orthologs' odds ratios correlations with temperature.
486    Table S6. Lineage, completeness and contamination of Metagenome Assembled Genomes as
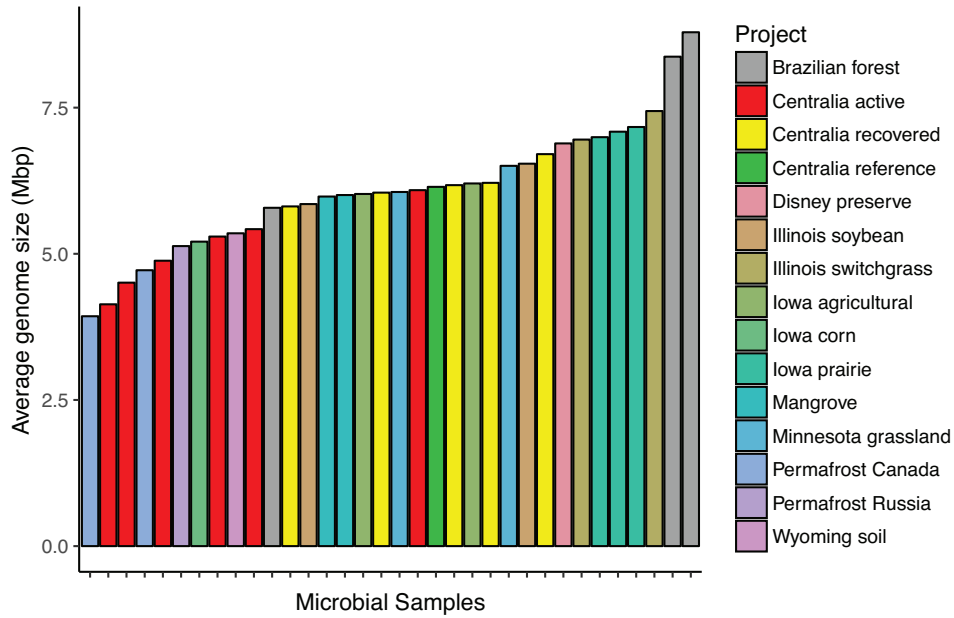487    estimated by CheckM
488

489
490    Figure 1. Changes in average genome and cell sizes across the soil temperature gradient in
491    Centralia. (A) Average genome size in each metagenome was calculated using MicrobeCensus
492    and plotted against site temperature.  (B) Average cell length was measured from 44-910 cells
493    from 3-9 replicate fields for each soil and plotted against soil temperature.  (C) Average genome
494    size had a direct relationship with average cell size.
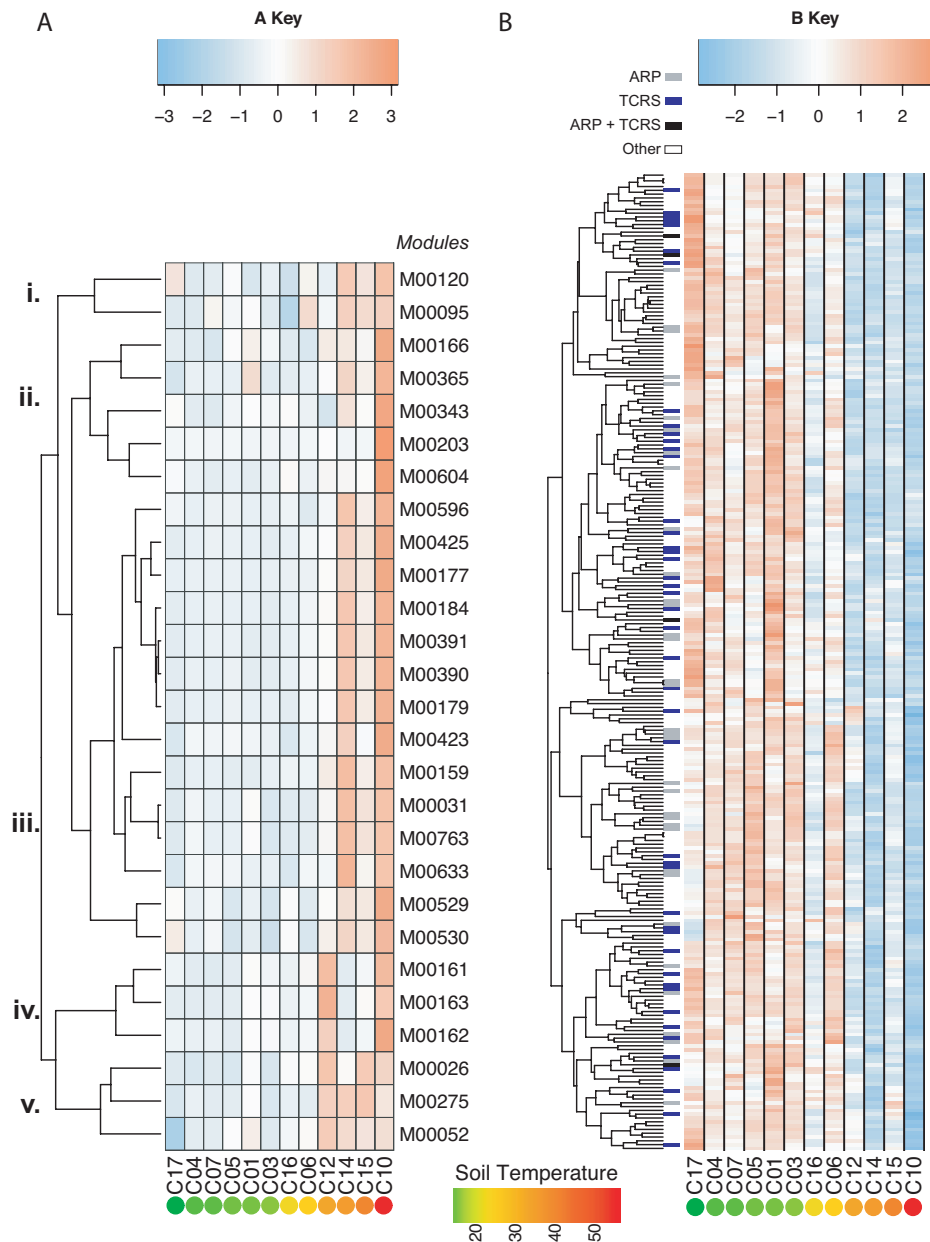495



496

497

498    Figure 2. Average genome sizes across soil metagenomes, estimated using MicrobeCensus.[18]
499    Samples are ordered by average genome size and colored by sample location.

500

501



502

503

504

505 Figure 3. Heatmap of KEGG modules correlated with temperature (false discovery rate adjusted
506 p-value < 0.05). Modules are centered and standardized (rows) across Centralia metagenomes
507 (columns), with warm colors showing relative enrichment and cool colors showing relative
508 depletion. Modules with significant relationships with temperature are shown. Sites are
509 arranged by increasing temperature from left to right. (A) 27 KEGG modules were positively
510 correlated with temperature (Pearson's R range = 0.646 - 0.933). (B) 257 KEGG modules were
511 negatively correlated with temperature (Pearson's R range = -0.642 to -0.925). A third of the
512 KEGG modules negatively correlated with temperature were either two-component regulatory
513 systems (TCRS, blue dendrogram tips), antimicrobial resistance or production (ARP, gray tips),
514 or both (black tips). Note differences in color gradient ranges across panels A and B.
515



516

517  Figure 4. KEGG modules that had notable enrichments or depletions with temperature. (A) The
518  median abundances of KEGG modules for denitrification (red), dissimilatory nitrate reduction
519  (green) and dissimilatory sulfate reduction (blue) were all positively correlated with
520  temperature. (B) Pearson's correlation values for all detected modules classified as antibiotic
521  resistance and production (gray density, n = 62 detected modules) or two-component
522  regulatory systems (blue density, n = 89 detected modules).  The black vertical line
523  distinguishes correlation values that are significant at a false discovery rate adjusted p-value <
524  0.05 (left), and all of these had a strong and negatively relationship with temperature. In total,
525  there were thirty-nine antimicrobial resistance and production modules and forty-nine two-
526  component regulatory system modules that significantly decreased with temperature.
527
528

529