

1 **MiDAS 5: Global diversity of bacteria and archaea in anaerobic digesters**

2

3 **Authors:** Morten Kam Dahl Dueholm^{1,*}, Kasper Skytte Andersen¹, Anne-Kirstine C.
4 Petersen¹, Vibeke Rudkjøbing¹, Madalena Alves², Yadira Bajón-Fernández³, Damien
5 Batstone⁴, Caitlyn Butler⁵, Mercedes Cecilia Cruz⁶, Åsa Davidsson⁷, Leonardo
6 Erijman⁸, Christof Holliger⁹, Konrad Koch¹⁰, Norbert Kreuzinger¹¹, Changsoo Lee¹²,
7 Gerasimos Lyberatos¹³, Srikanth Mutnuri¹⁴, Vincent O'Flaherty¹⁵, Piotr Oleskowicz-
8 Popiel¹⁶, Dana Pokorna¹⁷, Veronica Rajal⁶, Michael Recktenwald¹⁸, Jorge Rodríguez¹⁹,
9 Pascal E. Saikaly²⁰, Nick Tooker⁵, Julia Vierheilig¹¹, Jo De Vrieze²¹, Christian
10 Wurzbacher¹⁰ & Per Halkjær Nielsen^{1,*}

11

12 **Affiliations:**

13 ¹Center for Microbial Communities, Department of Chemistry and Bioscience,
14 Aalborg University, Aalborg, Denmark. ²Centre of Biological Engineering, University
15 of Minho, Portugal. ³School of Water, Energy and Environment, Cranfield University.
16 Cranfield, United Kingdom. ⁴Australian Centre for Water and Environmental
17 Biotechnology (ACWEB), The University of Queensland, Australia. ⁵Department of
18 Civil and Environmental Engineering, University of Massachusetts Amherst, MA,
19 USA. ⁶National University of Salta, Salta, Argentina. ⁷Department of Chemical
20 Engineering, Lund University, Sweden. ⁸INGEBI-CONICET, University of Buenos
21 Aires, Argentina. ⁹Laboratory for Environmental Biotechnology, Ecole Polytechnique
22 Fédérale de Lausanne (EPFL), Switzerland. ¹⁰Chair of Urban Water Systems
23 Engineering, Technical University of Munich (TUM), Garching, Germany. ¹¹Institute
24 of Water Quality and Resource Management, TU Wien, Austria. ¹²Department of
25 Urban and Environmental Engineering & Graduate School of Carbon Neutrality,
26 Ulsan National Institute of Science and Technology (UNIST), South Korea. ¹³School
27 of Chemical Engineering, National Technical University of Athens, Greece. ¹⁴Applied
28 Environmental Biotechnology Laboratory, Birla Institute of Technology and Science
29 (BITS-Pilani), India. ¹⁵School of Biological and Chemical Sciences and Ryan
30 Institute, University of Galway, Ireland. ¹⁶Water Supply and Bioeconomy Division,
31 Faculty of Environmental Engineering and Energy, Poznan University of Technology,
32 Poland. ¹⁷Department of Water Technology and Environmental Engineering,
33 University of Chemistry and Technology Prague, Czech Republic. ¹⁸Research
34 Scientist at Kemira Oyj, Espoo R&D Center, Finland. ¹⁹Chemical Engineering
35 Department, Khalifa University, United Arab Emirates. ²⁰Environmental Science and
36 Engineering Program, Biological and Environmental Science and Engineering
37 Division, King Abdullah University of Science and Technology (KAUST), Saudi
38 Arabia. ²¹Center for Microbial Ecology and Technology (CMET), Ghent University,
39 Belgium.

40 *Correspondence to: Per Halkjær Nielsen, Center for Microbial Communities,
41 Department of Chemistry and Bioscience, Aalborg University, Fredrik Bajers Vej 7H,
42 9220 Aalborg, Denmark; Phone: +45 2173 5089; Fax: Not available; E-mail:
43 phn@bio.aau.dk or Morten Kam Dahl Dueholm, Center for Microbial Communities,
44 Department of Chemistry and Bioscience, Aalborg University, Fredrik Bajers Vej 7H,
45 9220 Aalborg, Denmark; Phone: (+45) 9940 8508; Fax: Not available; E-mail:
46 md@bio.aau.dk

47

48 **Running title:** Global microbiota of anaerobic digesters

49 **Abstract**

50 Anaerobic digestion represents a key biotechnology for the transformation of organic
51 waste into renewable energy (biogas) and relies on complex microbial communities
52 that work in concert to degrade the complex substrates into methane and carbon
53 dioxide. Here, we sequenced more than half a million high-quality, full-length 16S
54 rRNA gene sequences from 285 full-scale anaerobic digesters (ADs) across the world
55 to expand our knowledge about diversity and function of the bacteria and archaea in
56 ADs. The sequences were processed into full-length 16S rRNA amplicon sequence
57 variants (FL-ASVs), which were added to the MiDAS 4 database for bacteria and
58 archaea in wastewater treatment systems to create MiDAS 5. The expansion of the
59 MiDAS database significantly increased the coverage for bacteria and archaea in ADs
60 worldwide, leading to an improved rate of genus and species-level classification. Using
61 MiDAS 5, we carried out an amplicon-based, global-scale microbial community
62 profiling of the sampled ADs using three common sets of primers targeting different
63 regions of the 16S rRNA gene in bacteria and/or archaea. We revealed how
64 environmental conditions and biogeography shape the AD microbiota. We also identify
65 core and conditionally rare or abundant taxa, encompassing 692 genera and 1013
66 species. These represent 84-99% and 18-61% of the accumulated read abundance
67 respectively, across samples depending on the amplicon primers used. Finally, we
68 examined the global diversity of functional groups with known importance for the
69 anaerobic digestion process. Our online global MiDAS Field Guide presents the data
70 generated in this study and summarizes present knowledge about all taxa.
71

72 **Introduction:**

73 Anaerobic digestion has gained attention as an important, sustainable biotechnology as
74 it provides several benefits that align with the goals of sustainability. It can help to
75 produce renewable energy (biogas) from organic waste such as manure, food waste,
76 and sludge from wastewater treatment plants (WWTPs) ^{1,2}. The anaerobic digestion
77 process also reduces pathogens and the amount of organic waste that is sent to landfills,
78 thereby reducing methane emissions and supporting sustainable waste management
79 practices ¹. Finally, the fertilizer that is produced as a byproduct of anaerobic digestion
80 can be used to support sustainable agriculture, reducing the need for synthetic fertilizers
81 that can have negative environmental impacts ^{3,4}.

82
83 The anaerobic digestion process relies on the microbial degradation and conversion of
84 organic matter, which requires a complex interplay between several functional guilds.
85 These include hydrolyzing, acidogenic, and acetogenic syntrophic bacteria as well as
86 methanogenic archaea ⁵. The taxonomy is poorly characterized for many of the
87 microorganisms in anaerobic digesters (ADs), and even among the most abundant taxa
88 many lack genus- or species-level classifications ⁶. To optimize performance, a
89 comprehensive knowledge about microbial immigration/competition,
90 environmental/operational conditions, and taxonomy is essential ⁷⁻⁹. Recent microbial
91 surveys have increased our knowledge about the anaerobic digestion process ^{7,10-16}.
92 However, sharing knowledge across studies is still hindered by the absence of
93 standardized protocols and a common reference database with a unifying taxonomy
94 ^{17,18}. To facilitate collaboration and knowledge sharing, it is essential to establish these
95 standard protocols and resources.

96
97 The Microbial Database for Activated Sludge and Anaerobic Digesters (MiDAS)
98 project was established as an open-source platform for sharing updated knowledge
99 about the physiology and ecology of the important microorganisms present in
100 engineered ecosystems of activated sludge plants, ADs, and related WWTPs ¹⁷⁻²⁰.
101 MiDAS provides standardized protocols for microbial profiling of microbes in
102 wastewater treatment systems ²¹, an ecosystem-specific full-length 16S rRNA gene
103 reference database ^{18,20}, and a field guide where knowledge about the specific genera
104 are stored and shared (<https://www.midasfieldguide.org>).

105
106 The MiDAS 16S rRNA gene reference database was created based on millions of high-
107 quality, chimera-free, full-length 16S rRNA genes resolved into amplicon sequence
108 variants (ASVs) and classified using automated taxonomy assignment (AutoTax) ^{6,18,20}.
109 AutoTax provides a comprehensive seven-rank taxonomy (kingdom to species-level)
110 for all reference sequences based on the most recent version of the SILVA SSURef 99
111 NR taxonomy and includes a robust placeholder taxonomy for lineages without an
112 official taxonomy ⁶. The placeholder taxa are easily distinguishable by their names,
113 formatted as 'midas_x_y', where 'x' indicates the taxonomic rank and 'y' is a numerical

114 identifier. This naming convention facilitates the study of unclassified alongside
115 classified taxa across various taxonomic ranks. The placeholder taxonomy should not
116 be seen as a replacement for proper taxonomic classifications but can pinpoint
117 important lineages that should be studied in depth using phylogenomics^{22–26}.

118

119 The MiDAS 16S rRNA gene reference database (MiDAS 4.8.1) currently contains
120 reference sequences from anaerobic digesters located at WWTPs in Denmark and
121 WWTPs worldwide²⁰. However, it may not provide comprehensive coverage for all
122 important microbes found in ADs treating other types of waste or in other locations.

123

124 In this study, we introduce MiDAS 5, an updated version of MiDAS 4 expanded with
125 more than half a million high-quality, full-length archaeal and bacterial 16S rRNA gene
126 sequences from 285 ADs worldwide treating different types of biowaste. We carried
127 out a global survey of ADs using three commonly used short-read amplicon primer sets
128 targeting bacteria (V1-V3), archaea (V3-V5), and both (V4). This data was used in
129 combination with MiDAS 5 to i) link the global diversity of bacteria and archaea to
130 biogeography and environmental factors, ii) identify important core taxa, and iii)
131 uncover the global diversity within selected functional guilds. The results provide a
132 solid foundation for future research on AD microbiology.

133

134 **Results and Discussion:**

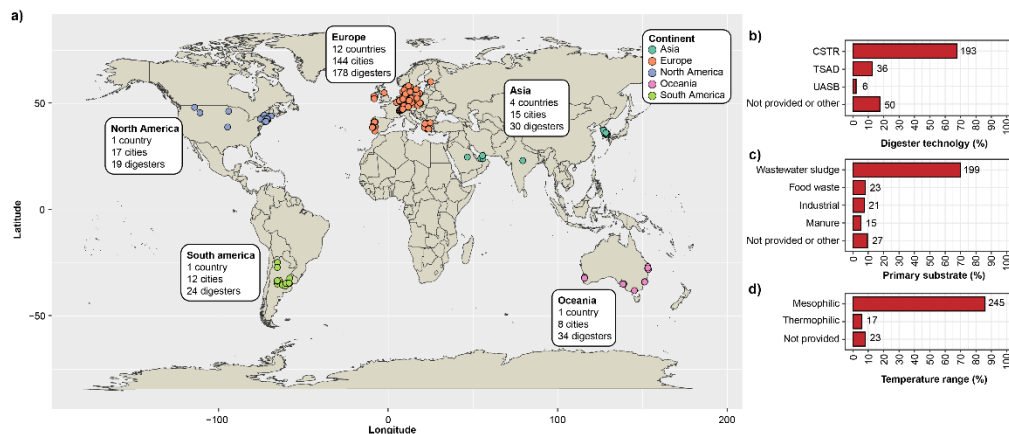
135 The MiDAS Global Consortium for Anaerobic Digesters was established in 2018 to
136 coordinate the sampling and collection of metadata from ADs worldwide
137 (Supplementary Data 1). Samples were obtained in duplicates from 285 ADs in 196
138 cities in 19 countries on five continents (Fig. 1a). Most of the ADs treated surplus
139 sludge from wastewater treatment plants (69.8%) (Fig. 1b). However, ADs treating
140 food waste (8.1%), industrial waste (7.4%), and manure (5.3%) were also included in
141 the survey. Most of the ADs were mesophilic (86.0%), few were thermophilic (6.0%),
142 and some missed temperature data (8.1%). The main digester technology used was
143 completely mixed reactors (67.7%) followed by two-stage reactors (12.6%). A few
144 upflow anaerobic sludge blanket (UASB) and other types were also sampled to expand
145 the diversity of digester types.

146

147 ***Expanding the MiDAS database with reference from global ADs***

148 To expand the MiDAS database with sequences from ADs across the globe, we applied
149 high-fidelity, full-length 16S rRNA gene sequencing on all samples collected in this
150 study. More than half a million full-length 16S rRNA gene sequence reads, representing
151 both bacteria and archaea, were obtained after quality filtering and primer trimming.
152 After processing the sequence reads with AutoTax to produce full-length 16S rRNA
153 gene ASVs (FL-ASVs), these were compared and added to the existing 90,164 FL-
154 ASVs in the MiDAS 4.8.1 database. The combined number was then deduplicated,
155 resulting in a total of 120,408 non-redundant FL-ASV reference sequences in the

156 expanded MiDAS 5 database. This represents an increase of 30,246 new FL-ASVs
 157 when compared to the previous version.
 158



159
 160 Fig. 1: Sampling of anaerobic digesters (ADs) across the world. a) Geographical distribution of ADs
 161 included. b) Distribution of digester technologies. CSTR: Continuous stirred-tank reactor; TSAD: Two-
 162 stage anaerobic digestion; UASB: Upflow anaerobic sludge blanket. c) Distribution of primary
 163 substrates. d) Distribution of digester temperatures. The values next to the bars are the number of ADs
 164 in each group.
 165

166 Table 1: Sequence novelty of FL-ASVs obtained in this study. Sequence novelty was determined based
 167 on the percentage identity between each of the 30,246 new FL-ASV and their closest relative in the
 168 databases indicated and identity thresholds for each taxonomic rank proposed by Yarza *et al.*²⁷ shown
 169 in the parentheses.

	SILVA 138.1 SSURef NR99		MiDAS 4.8.1	
	Sequences	Percentage	Sequences	Percentage
Phylum (<75.0%)	28	0.09%	67	0.22%
Class (<78.5%)	55	0.18%	186	0.61%
Order (<82.0%)	112	0.37%	462	1.53%
Family (<86.5%)	354	1.17%	1,483	4.90%
Genus (<94.5%)	5,240	17.32%	9,419	31.14%
Species (<98.7%)	15,806	52.26%	16,863	55.75%

170
 171 The novelty of the 30,246 new FL-ASVs were determined based on the percent identity
 172 shared with their closest relatives in the SILVA 138.1 SSURef NR99 and MiDAS 4.8.1
 173 database using the threshold for each taxonomic rank proposed by Yarza *et al.*²⁷ (Table
 174 1). It should be noted that these thresholds do not uniformly apply across the bacterial
 175 phylogenetic tree; therefore, our taxonomic assignments should be considered as
 176 approximations intended to facilitate biological interpretation. 17% and 31% of the new
 177 FL-ASV lacked genus-level homologs ($\geq 94.5\%$ identity) and 52% and 56% were
 178 without species-level homologs ($\geq 98.7\%$ identity) in SILVA 138.1 and MiDAS 4,
 179 respectively. This suggests a substantial increase in the diversity within the MiDAS 5
 180 database.

181 ***MiDAS 5 introduces many new taxa***

182 To investigate how the new FL-ASVs affected the taxonomic diversity in the MiDAS
183 database, we determined the number of additional taxa introduced at different
184 taxonomic ranks (Table 2). A substantial increase in diversity was observed with the
185 addition of 2,770 new genera (29.2% increase) and 8,858 new species (28.3% increase).
186 However, many additional taxa were also introduced at higher taxonomic ranks
187 including six more bacterial and five more archaeal phyla previously known from the
188 SILVA taxonomy. In addition, we identified nine lineages classified as MiDAS
189 placeholder phyla. However, phylogenetic analysis revealed that these lineages branch
190 closely to mitochondrial sequences, indicating they are likely mitochondrial in origin.
191 The largest percentage of the new FL-ASVs (42.8%) were found within the Firmicutes
192 (Supplementary Fig. 1a). Firmicutes often occur in high abundance in ADs, where they
193 are involved in fermentation and thereby directly stimulate biogas yields^{7,10,13,15,28}. A
194 closer look into the expanded diversity within the Firmicutes revealed that new FL-
195 ASVs were associated with several families (Supplementary Fig. 1b), including
196 Hungateiclostridiaceae (1,324 FL-ASVs), Lachnospiraceae (788 FL-ASVs),
197 Peptostreptococcales-Tissierellales Family_XI (763 FL-ASVs), Christensenellaceae
198 (754 FL-ASVs), Caldicoprobacteraceae (620 FL-ASVs), and Syntrophomonadaceae
199 (555 FL-ASVs). The Syntrophomonadaceae are of special relevance, as this family
200 includes several syntrophic fatty acid degrading bacteria, which are often the metabolic
201 bottleneck in the overall ADs process^{29,30}.

202 Table 2: New taxa introduced with MiDAS 5. The number of new taxa represent unique taxa at the
203 different taxonomic ranks that were not part of MiDAS 4.8.1²⁰ and includes both official taxonomic
204 names and de novo placeholder names provided by AutoTax⁶.

	Total taxa	New taxa	Increase (%)
Phylum	105	20	23.5%
Class	259	37	16.7%
Order	727	104	16.7%
Family	2,212	360	19.4%
Genus	12,254	2,770	29.2%
Species	40,207	8,858	28.3%

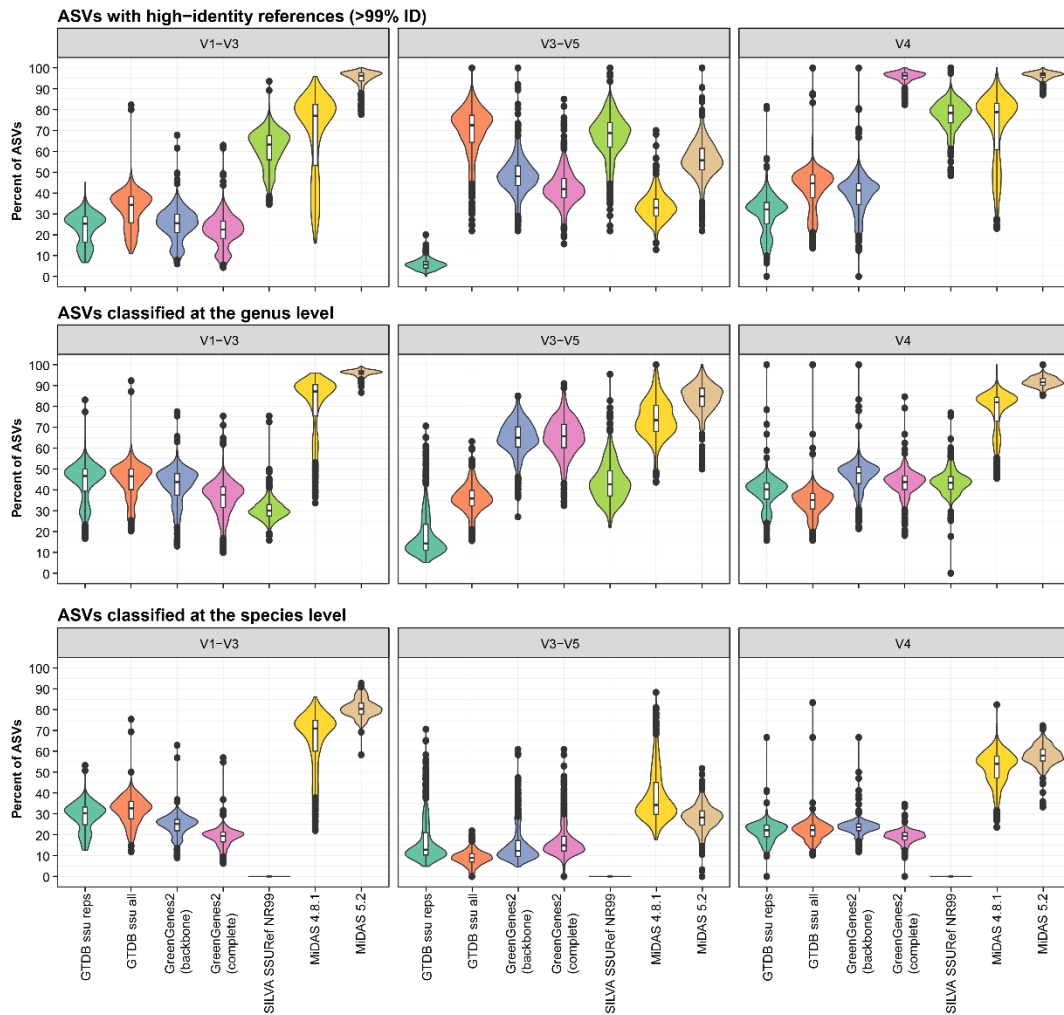
205

206 ***MiDAS 5 provides improved coverage and classifications for AD microbiota***

207 The performance of the MiDAS 5 database was evaluated based on three ASV-
208 resolved, short-read, 16S rRNA gene amplicon datasets generated from the AD samples
209 collected in this study (Fig. 2). The V1-V3 amplicons include only bacteria and provide
210 high phylogenetic resolution. However, the primers targeting this region have a lower
211 coverage for the known bacterial diversity according to in silico evaluations^{6,31}. The
212 V4 amplicons include both bacteria and some archaeal lineages and are commonly used
213 due to a very good coverage of the known bacterial diversity. However, the amplicons
214 have a weaker phylogenetic resolution compared to V1-V3, which in many cases
215 prevent species-level classifications^{6,31}. The V3-V5 amplicons cover mainly archaea
216 and have previously been used to describe their diversity in ADs^{7,10}.

217
218 Our initial analysis involved non-heuristic mapping of short-read ASVs against MiDAS
219 5 and other widely used reference databases, including the newly released GreenGenes2
220 ³². This step allowed us to establish the percent identity between each ASV and its
221 closest match across the databases. We then calculated the percentage of ASVs that
222 have high-identity matches ($\geq 99\%$ identity) in each sample and database. To focus on
223 active microbial populations, we excluded ASVs representing the rare biosphere (those
224 with $< 0.01\%$ relative abundance), which are often enriched in non-growing organisms
225 and environmental DNA^{7,10}. MiDAS 5 performed exceptionally well for bacteria with
226 high-identity hits of $94.8\% \pm 4.2\%$ (mean \pm SD) for V1-V3 and $96.3\% \pm 2.1\%$ for V4
227 ASVs, compared to $67.9\% \pm 19.7\%$ and $71.4\% \pm 16.1\%$ for MiDAS 4, and $61.1\% \pm 9.2\%$
228 and $77.1\% \pm 7.8\%$ for SILVA v.138.1 (Fig. 2). The complete GreenGenes2 database
229 displayed a coverage close to that of MiDAS 5 for V4 ASVs ($95.4\% \pm 3.3\%$) but a much
230 lower coverage for V1-V3 ($32.1\% \pm 8.9\%$). The reason is that the complete
231 GreenGenes2 database contains V4 ASVs from Qiita³³ in addition to full-length 16S
232 rRNA gene sequences³². For the V3-V5 archaeal dataset, an increase in coverage was
233 observed from $33.5\% \pm 7.0\%$ with MiDAS 4 to $55.9\% \pm 9.5\%$ with MiDAS 5. However,
234 the SILVA database ($67.0\% \pm 11.0\%$) and the complete GTDB database ($69.2\% \pm 13.0\%$)
235 provide even better coverage. The lower coverage for archaea compared to bacteria in
236 MiDAS 5 is likely due to reduced sequencing efforts and the challenges in designing
237 effective universal primers for archaeal full-length 16S rRNA gene sequencing^{34,35}.

238
239 Because the sampling of ADs was directed towards mesophilic digesters treating
240 surplus sludge from WWTPs, we also evaluated the MiDAS 5 coverage for ADs
241 treating different primary substrates and temperatures (Supplementary Fig. 2). MiDAS
242 5 gave very good coverage for all sample types supporting the general applicability of
243 the reference database for ADs. Finally, to provide additional support for the general
244 applicability of the MiDAS 5 database, we evaluated it based on previously published
245 V4-V5 amplicon data from 90 full-scale ADs at 51 municipal WWTPs unrelated to this
246 study¹⁴. MiDAS 5 contained high-identity hits for $91.8\% \pm 6.8\%$ of the ASVs
247 (Supplementary Fig. 3).



248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263

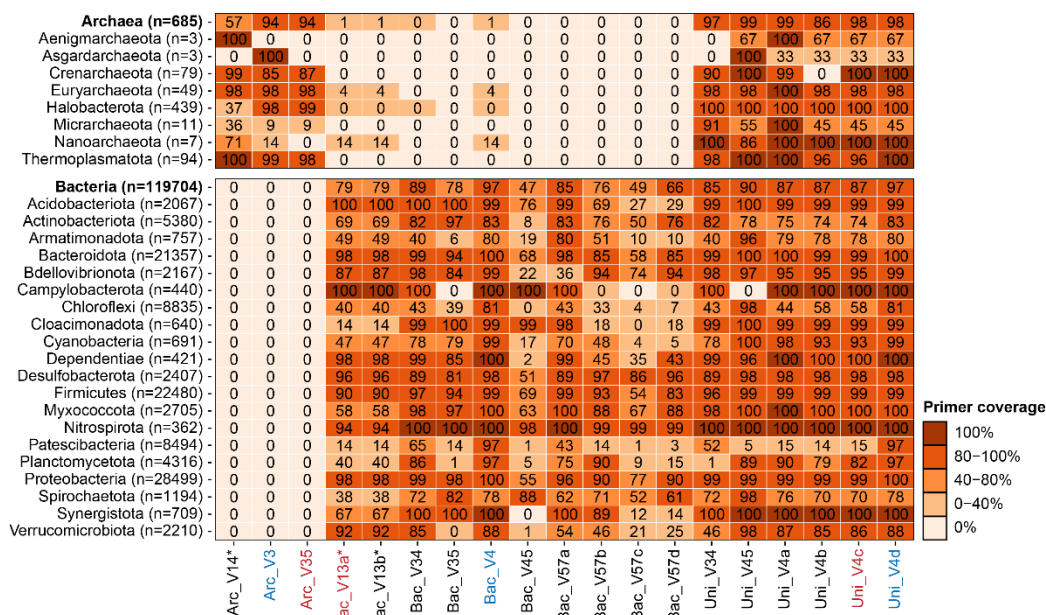
Fig. 2: Database evaluation based on short-amplicon data from this study. The ASVs for each of the samples were filtered based on their relative abundance (only ASVs with $\geq 0.01\%$ relative abundance were kept) before the analyses. The percentage of the microbial community represented by the remaining ASVs after the filtering was $95.44\% \pm 2.23\%$ (mean \pm SD) for V1-V3 amplicons (only bacteria), $99.65\% \pm 0.17\%$ for V3-V5 amplicons (mainly archaea), and $97.34\% \pm 2.01\%$ for V4 amplicons (bacteria and archaea) across samples. High-identity ($\geq 99\%$) hits were determined by stringent mapping of ASVs to each reference database. Classification of ASVs was done using the SINTAX classifier. The violin and box plots illustrate the distribution of the percentage of ASVs with high-identity hits or genus/species-level classifications for each database, analyzed across 570 biologically independent samples, including two biological replicates for each digester. Box plots indicate median (middle line), 25th, 75th percentile (box), and the min and max values after removing outliers based on 1.5x interquartile range (whiskers). Outliers have been removed from the box plots to ease visualization. Different colors are used to distinguish the different databases: GTDB_bac120_ssu_reps_r214, GTDB_ssu_all_r214, GreenGenes2_2022_10 (backbone and complete database), SILVA 138.1 SSURef NR99, MiDAS 4.8.1, and MiDAS 5.2.

264 Our second database evaluation was based on the classification of ASVs from each
 265 amplicon dataset using the SINTAX classifier (Fig. 2). We found that MiDAS 5 greatly
 266 improved the rates of genus-level classification ($96.3\% \pm 1.4\%$ for V1-V3, $91.5\% \pm 2.6\%$
 267 for V4, and $82.6\% \pm 7.5\%$ for V3-V5) compared to MiDAS 4 ($80.2\% \pm 14.9\%$ for V1-
 268 V3, $77.3\% \pm 10.5\%$ for V4, and $74.7\% \pm 9.3\%$ for V3-V5), and the rates of classification
 269 were more than two fold higher than those obtained with any of the other evaluated
 270 databases for bacteria and also higher for archaea. Analysis of species-level
 271 classifications revealed similar improvements with MiDAS 5 for bacteria (Fig. 2).
 272 However, a decrease in species-level classifications was observed between MiDAS 4
 273 and 5 for the archaeal V3-V5 dataset. We hypothesize that this effect relates to over-
 274 classifications with MiDAS 4 due to the lack of appropriate reference sequences in
 275 MiDAS 4.

276

277 Finally, we investigated if the additional reference sequences introduced in MiDAS 5
 278 could improve classification of amplicon data from WWTPs based on data from the
 279 MiDAS global sampling of WWTPs²⁰ and the Global Water Microbiome Consortium
 280 project³⁶ (Supplementary Fig. 4). Interestingly, no statistically significant
 281 improvements were observed. This highlights that most of the added references
 282 originated from anaerobic digester-specific taxa.

283



284

285 Fig. 3: Ecosystem-specific primer coverage for all archaea and bacteria and the 20 most diverse bacterial
 286 and all archaeal phyla based on unique FL-ASVs. The number of FL-ASVs for each taxon (n) is provided
 287 next to the names. The coverage was determined as the percentage of FL-ASVs in the MiDAS 5.2
 288 database with perfect hits for both forward and reverse primers. The primer pairs marked in red were
 289 used in the current study and the ones marked in blue are those recommended here based on coverage.
 290 Detailed information of all primer pairs and coverage information for all taxa in MiDAS 5.2 are provided
 291 in Supplementary Data 2. *Only the reverse primer was evaluated for these primer pairs because the
 292 forward primer was used to create the reference sequences in MiDAS. The coverage might therefore be
 293 overestimated for these primer pairs.

294 ***Evaluation of 16S rRNA gene amplicon primers for community profiling of ADs***

295 The comprehensive ASV-resolved MiDAS 5 database provides a unique opportunity
296 to determine the theoretical coverage of commonly applied 16S rRNA gene amplicon
297 primer pairs for bacteria and archaea in ADs (Fig. 3). This information is highly
298 valuable when designing experiments, especially if targeting specific taxa.
299 Accordingly, we have determined the theoretical coverage for several commonly
300 applied primer pairs for all kingdom to species-level taxa in MiDAS 5 (Supplementary
301 Data 2). We found a fairly low coverage of the V1-V3 primer pair (perfect hits for
302 $\leq 79\%$ of the bacterial FL-ASVs), which we commonly use due to its high phylogenetic
303 resolution^{6,20}. We should therefore expect a significant bias when using this primer
304 pair. The V4 primer used here and in the Earth Microbiome project³⁷ showed good
305 coverage for both bacteria (perfect hits for 87% of the FL-ASVs) and archaea (perfect
306 hits for 98% of the FL-ASVs). However, a recently published primer pair for the V4
307 region, designed to improve coverage for Patescibacteria³⁸, showed even better
308 coverage for bacteria, achieving perfect hits for 97% of the FL-ASVs. Although this
309 primer pair does not target archaea, adding degeneracy at a single base in one of the
310 primers also provided coverage for archaea, with perfect hits for 98% of the FL-ASVs.
311 The exceptional coverage offered by this new primer pair leads us to recommend it for
312 the profiling of anaerobic digesters (ADs), despite its lower phylogenetic signal
313 compared to the V1-V3 primers. The V3-V5 primer pair, which was used here to target
314 archaea only, also had good coverage for archaea, though not as extensive as that of the
315 V4 primers, supporting the choice of the latter.

316

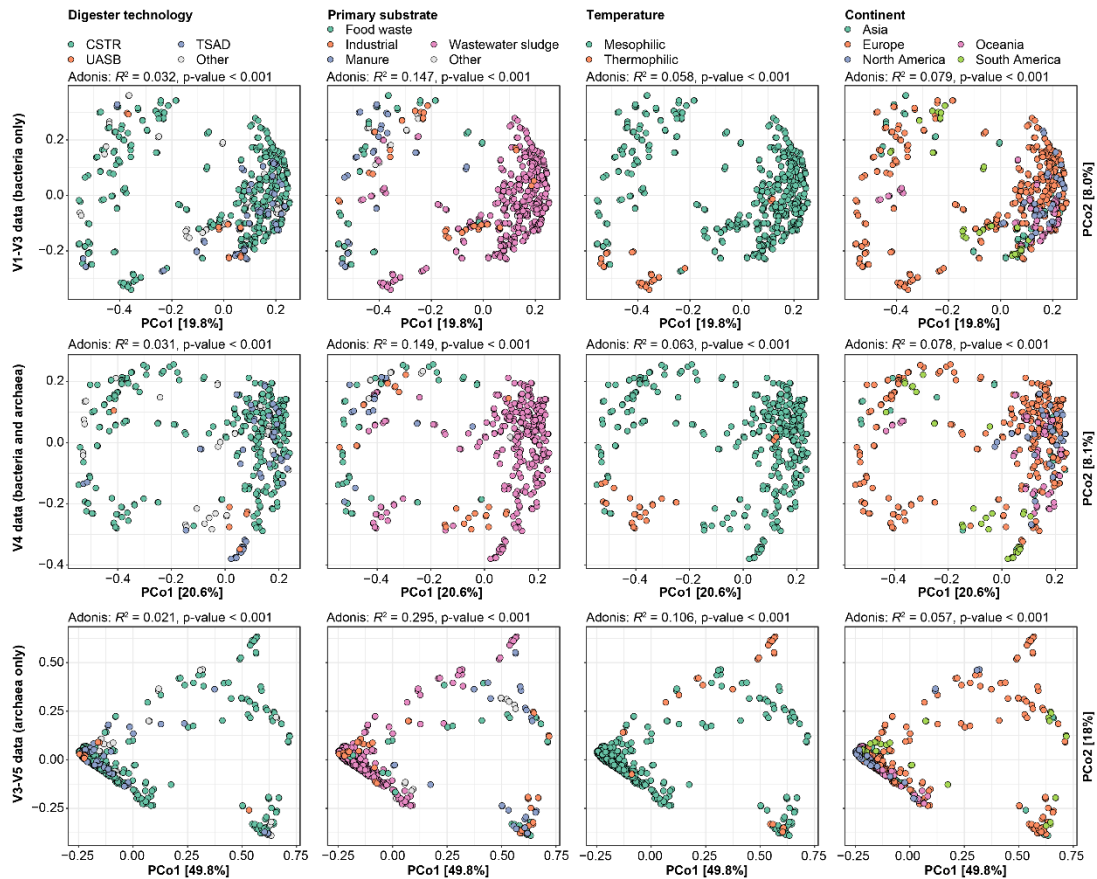
317 ***Effect of process and environmental factors on the AD microbiota***

318 Alpha diversity analyses showed that the rarefied (10,000 read per sample) ASV
319 richness and inverse Simpsons diversity in ADs were affected mainly by the primary
320 substrate type and the temperature in the ADs (Supplementary Fig. 5). Significantly
321 higher bacterial richness and diversity were observed for ADs treating surplus sludge
322 from WWTPs compared to the other types of substrates. This effect likely reflects the
323 extensive immigration of bacteria into the ADs with the surplus sludge^{7,10,39}. A higher
324 richness and diversity were observed for bacteria in mesophilic ADs compared to
325 thermophilic ADs. A similar trend has previously been observed for full-scale ADs
326 treating manure^{40,41}, household waste⁴², and surplus sludge from WWTPs⁷.

327

328 Genus-level taxonomic beta-diversity was used to investigate the effect of process
329 conditions and geography on the overall microbiota in ADs using principal coordinate
330 analysis (PCoA) and permutational multivariate analysis of variance (PERMANOVA)
331 (Fig. 4). We used this approach because many of the important traits are categorical
332 (yes/no) and only conserved at lower taxonomic ranks (genus/species)⁴³. Furthermore,
333 MiDAS 5 enabled us to classify almost all our ASVs at the genus-level, thereby
334 providing a comprehensive description of the microbiota. The PERMANOVA (Adonis
335 R^2 values) showed that the overall microbial community was mainly explained by the

336 primary substrate and to a lesser extent by temperature, continent, and digester
 337 technology (Fig. 4). This trend was observed for both bacteria and archaea. The
 338 percentage of total variation explained by each parameter was, except for the primary
 339 substrate, low, suggesting that the global AD microbiota represents a continuous
 340 distribution rather than distinct states, as also observed for the human gut microbiota⁴⁴
 341 and WWTPs²⁰. The pronounced effect of the primary substrates highlights that these
 342 substrates are distinct and abundant in microbes, particularly in the case of manure and
 343 wastewater sludge.
 344



345
 346 Fig. 4: Effects of process and environmental factors on the anaerobic digester microbiota. Principal
 347 coordinate analyses of Bray-Curtis beta-diversity for genera based on V1-V3 (bacteria), V4 (archaea and
 348 bacteria), and V3-V5 (archaea) amplicon data. Samples are colored based on metadata. The fraction of
 349 variation in the microbial community explained by each variable in isolation was determined by
 350 PERMANOVA (Adonis R^2 -values). Exact p-values less than 0.001 could not be confidently determined
 351 due to the chosen number of permutations. CSTR: Continuous stirred-tank reactor; TSAD: Two-stage
 352 anaerobic digestion; UASB: Upflow anaerobic sludge blanket.

353 ***Core and conditional rare or abundant taxa in the global AD microbiota***

354 The global AD microbiota represents a huge microbial diversity. However, most
355 organisms only occur in very low abundance and are therefore unlikely to have any
356 quantitative impact on the overall metabolism and the process performance in ADs.
357 Analysis of core and conditionally rare or abundant taxa (CRAT) is a powerful
358 approach to identify the most important genera and species within a specific ecosystem
359 ^{20,28,45}. The CRAT may include taxa related to process disturbances, such as filamentous
360 microbes associated with foam formation, or taxa associated with the degradation of
361 special substrates found in, e.g., industrial waste.

362

363 We recently introduced and applied the following core and CRAT definitions in our
364 survey of the global microbiota of wastewater treatment plants: strict core (>0.1%
365 relative abundance in >80% of samples), general core (>0.1% relative abundance in
366 >50% of samples), loose core (>0.1% relative abundance in >20% of samples), and
367 CRAT (not part of the core, but present in at least one sample with a relative abundance
368 >1%) ²⁰. Here, we applied the same criteria to identify core and CRAT genera and
369 species in our global AD dataset. Because the primary substrate showed a strong effect
370 on the overall microbial community (Fig. 4), we determined the core and CRAT for
371 each individual substrate separately (Supplementary Data 3). Only mesophilic ADs
372 were examined for ADs treating food waste, industrial waste, and manure due to the
373 low number of thermophilic ADs sampled. Both mesophilic and thermophilic digesters
374 were examined for ADs treating wastewater sludge. To minimize the impact of primer
375 bias, we analyzed all three amplicon datasets and combined the results, including all
376 core and CRAT that were found in at least one of the datasets.

377

378 The core analysis revealed that most core genera were uniquely associated with specific
379 primary substrates and temperature range (Fig. 5a). However, there were also a
380 significant number of core genera shared across substrates (Fig. 5a). In contrast, very
381 few core species were shared between ADs treating different primary substrates (Fig.
382 5b). This fits well with similar results from a study of ADs in Belgium and Luxemburg
383 ¹³. To define a 'most wanted' list for bacteria and archaea in ADs globally, we linked
384 each core and CRAT to their highest-ranking category across primary substrates,
385 process temperatures, and primer pair (Supplementary Data 3). The resulting list
386 contained 501 core genera (75 strict, 117 general, and 309 loose) and 191 CRAT genera.
387 The strict core genera included 11 known methanogens and four known syntrophs (*Ca.*
388 *Phosphitivorax*, *Smithella*, *Syntrophomonas*, *Syntrophorhabdus*). At the species-level,
389 we identified 565 core (29 strict, 126 general, and 410 loose) and 448 CRAT species.
390 The strict core species included two methanogens (*Methanobrevibacter smithii* and
391 *Methanothermobacter midas_s_3958*) and one syntroph (*Syntrophomonas*
392 *midas_s_90707*). It is worth noting that a large fraction of the taxa observed in
393 ADs does not grow in the digesters, but only occurs because the taxa are in high
394 abundance in the feed ^{7,10,39}. Previous published data from Danish ADs treating

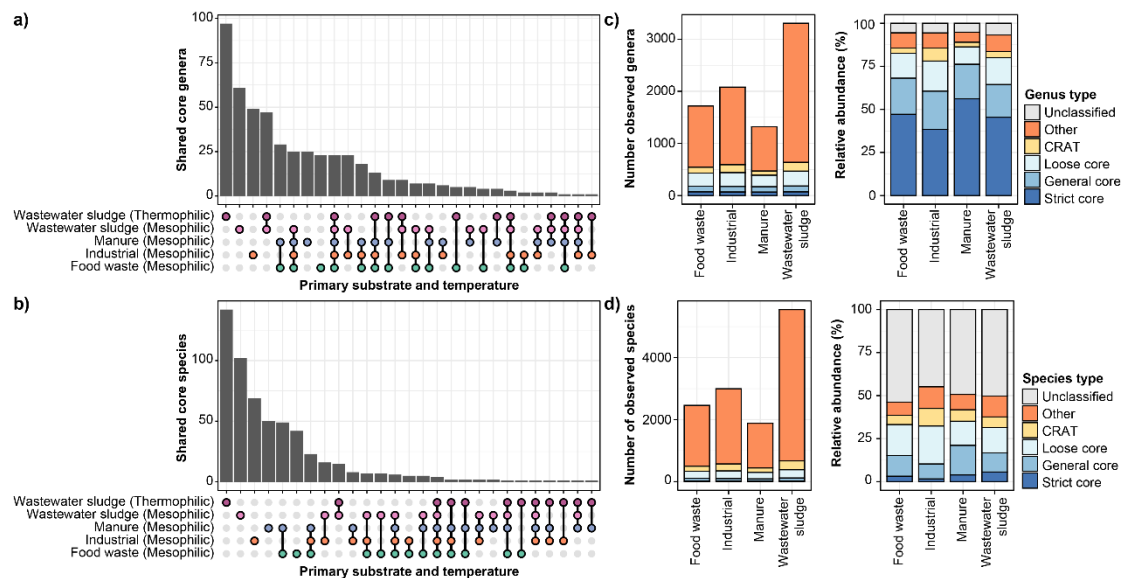
395 wastewater sludge⁷ classified 45 (9.0%) of the core genera observed in this study as
 396 non-growing (<20% of ASVs belonging to the specific taxa were classified as
 397 growing), whereas 393 (78.4%) were classified as growing. A similar analysis of core
 398 species classified 45 (8.0%) as non-growing and 391 (69.2%) as growing. However, it
 399 remains to be determined if these numbers also translate to global ADs.

400

401 **Many core and CRAT represent MiDAS placeholder taxa**

402 A large proportion of core and CRAT identified was classified as MiDAS *de novo* taxa.
 403 At the genus-level, 272/501 (54%) of the core genera and 119/191 (62%) of the CRAT
 404 genera had only MiDAS placeholder names, and at the species-level, the proportion
 405 was even higher. Here placeholder names were assigned to 514/565 (91%) of the core
 406 species and 422/448 (94%) CRAT species. These proportions are similar to those
 407 observed for the global microbiota in activated sludge²⁰ and reveals the importance of
 408 a taxonomic framework that can handle uncultured taxa which have not yet been
 409 officially classified.

410



411

412 Fig. 5: Core and conditionally rare or abundant taxa (CRAT) in anaerobic digesters globally. a) and b)
 413 UpSet plots displaying the number of shared core genera and species, respectively, across ADs treating
 414 different primary substrates and operating at different temperatures. c) and d) Number of observed genera
 415 and species, respectively, and their abundance in mesophilic ADs treating different primary substrates
 416 based on V4 amplicon data (bacteria and archaea). Values for genera and species are divided into strict
 417 core, general core, loose core, CRAT, other taxa, and unclassified ASVs based on the most wanted list
 418 (Supplementary Data 3). The relative abundance of different groups was calculated based on the mean
 419 relative abundance of individual genera or species across samples. Similar figures for V1-V3 (bacteria
 420 only) and V3-V5 (archaea only) amplicons data can be found in Supplementary Fig. 6.

421 ***The global AD microbiota is dominated by core and CRAT taxa***

422 Despite only accounting for a minor fraction of the total diversity in the ADs examined,
423 the core and CRAT represented most of the microbes ~~in the ecosystem~~ according to
424 relative amplicon read abundance (Fig. 5c, 5d, Supplementary Fig. 6). The core and
425 CRAT genera accounted for 85-92% (V1-V3), 84-89% (V4), and 96-99% (V3-V5) of
426 the accumulated read abundance in mesophilic ADs depending on primary substrates.
427 The remaining fractions consisted mainly of ASVs unclassified at the genus level, and
428 genera present in very low abundance, presumably with minor importance for the AD
429 performance.

430

431 For the species level, the core and CRAT represented 53-61% (V1-V3), 38-43% (V4),
432 and 18-47% (V3-V5) accumulated read abundance depending on the primary substrate.
433 The remaining fractions were mainly composed of ASVs, which could not be classified
434 at the species level, probably due to insufficient phylogenetic resolution in the short-
435 read amplicons^{6,31}. The lack of species-level classification was especially pronounced
436 for the archaeal V3-V5 ASVs in ADs treating industrial waste, manure, and wastewater
437 sludge (Supplementary Fig. 6).

438

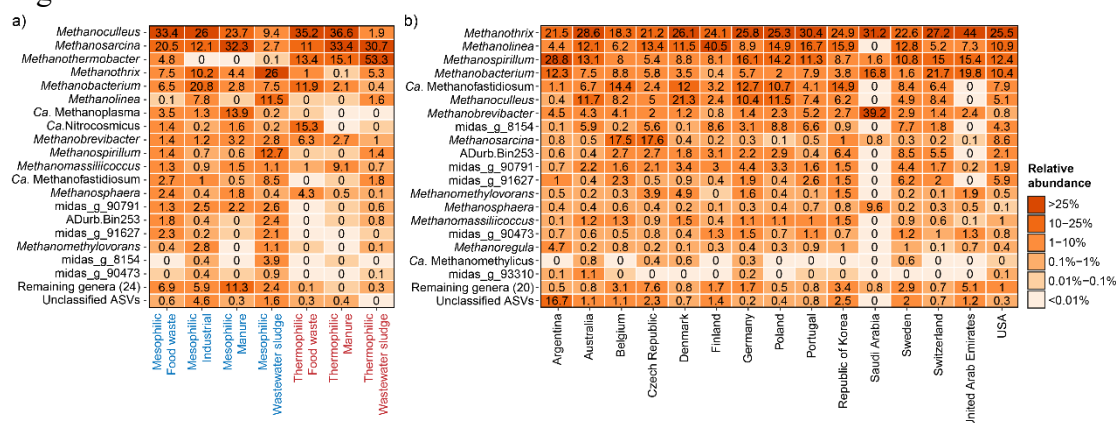
439 The large relative abundance of core and CRAT in the global AD microbiota suggests
440 that we can explain a substantial part of the metabolic processes in ADs, if we
441 understand the physiology and metabolic potential of these taxa.

442

443 ***Global diversity of archaea reveals new potential methanogens***

444 As methanogenic archaea are ultimately responsible for the generation of methane in
445 ADs, we examined the global diversity of archaea in all samples based on V4 (Fig. 6)
446 and V3-V5 amplicon data (Supplementary Fig. 7). The V4 amplicon data,
447 encompassing both archaea and bacteria, showed that the archaeal reads constituted
448 $5.6\% \pm 4.4\%$ for ADs treating food waste, $6.8\% \pm 4.4\%$ for manure, $6.4\% \pm 2.5\%$ for
449 wastewater sludge, and $13.7\% \pm 11.1\%$ for industrial waste. Many of the abundant
450 archaea represented well-known methanogens. However, we also observed several
451 abundant genera, only classified based on the MiDAS placeholder taxonomy, affiliated
452 to orders and families of known methanogens. These include *midas_g_91627* and
453 *midas_g_8154*, which represent new families within the orders Methanomicrobiales
454 and Methanofastidiosales, respectively, and *midas_g_90473* and *midas_g_93310*,
455 representing new genera within Methanomassiliicoccaceae and Methanospirillaceae,
456 respectively. In addition, we observed two abundant MiDAS placeholder genera
457 (*midas_g_90791* and *midas_g_97217*) that represent a new order within the class *Ca.*
458 *Bathyarchaeia*. Members of this class can have a versatile metabolism, and some
459 encode the key methanogenic enzyme methyl-coenzyme M reductase (MCR)^{46,47}.
460 Targeted metagenomics and assembly of metagenome-assembled genomes (MAGs)
461 should be applied to confirm the methanogenic potential of these new potential

462 methanogens, and our amplicon datasets provide insight into where these taxa occur in
 463 high abundance.



464 Fig. 6: Top 25 archaeal genera based on V4 amplicon data. The percent relative abundance represents
 465 the mean abundance relative to all archaea across a) different temperature range and primary substrates,
 466 and b) different countries considering only mesophilic ADs treating mainly wastewater sludge.
 467

468 The methanogenic community composition was clearly affected by the primary
 469 substrate and whether the digestion was performed at mesophilic or thermophilic
 470 conditions (Fig. 6a, Supplementary Fig. 7a). The most common methanogens across
 471 substrates and temperatures were *Methanoculleus*, *Methanosarcina*,
 472 *Methanothermobacter*, and *Methanotrix*. *Methanothermobacter* was as expected most
 473 abundant in thermophilic ADs. However, to our surprise, it also occurred in high
 474 relative abundance in several mesophilic reactors treating mainly food waste. We were
 475 not able to explain their occurrences in these ADs based on the available metadata for
 476 the plants, but future studies might shed light on the underlying mechanisms or
 477 environmental factors that enable this unexpected distribution.
 478

479 Because most of our samples originated from mesophilic reactors treating wastewater
 480 sludge, we examined the diversity of methanogens across countries in these ADs (Fig.
 481 6b, Supplementary Fig. 7b). This analysis revealed that the same genera were
 482 dominating across the world. The most common methanogens in these ADs were
 483 *Methanotrix*, *Methanolinea*, *Methanospirillum*, *Methanobacterium*, and the recently
 484 discovered *Ca. Methanofastidiosum*⁴⁸. Next, we examined if the methanogens were
 485 also conserved at higher phylogenetic resolution. As many archaeal ASVs could not be
 486 classified at the species-level, we examined the global diversity at the ASV-level
 487 (Supplementary Fig. 8). We found that the vast majority of the abundant ASVs occurred
 488 globally. The significant similarity of methanogens across various regions indicates
 489 substantial potential for global knowledge transfer concerning their management and
 490 utilization.
 491

492 Among the highly abundant archaea, we also observed an ammonia oxidizing archaeon
 493 (AOA) from the genus *Ca. Nitrosocosmicus*⁴⁹, which was especially abundant in
 494 thermophilic ADs treating food waste. This is surprising and may indicate that they also
 495

496 have an anaerobic physiology which should be investigated further. Another abundant
 497 archaeon was the *Ca. Diapherotrites ADurb.bin253* belonging to the order
 498 Woearchaeales which are characterized by ultra-small genomes and an anaerobic and
 499 parasitic/fermentation-based lifestyle⁵⁰.

500

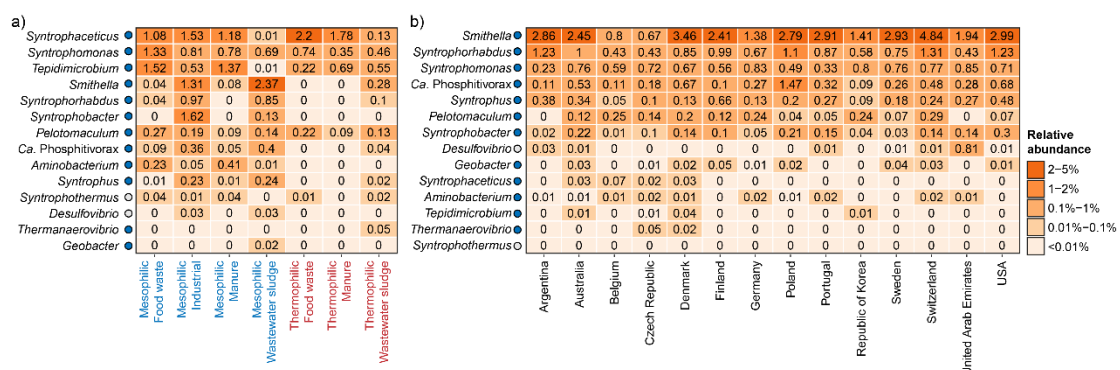
501 *Global diversity of syntrophic bacteria*

502 Syntrophic bacteria play a vital role in ADs by converting substrates, such as short-
 503 chain fatty acids, into acetate, H₂, and formate^{29,51,52}. These compounds serve as
 504 substrates or reducing equivalents for methanogens, which in turn produce methane and
 505 CO₂. This obligately mutualistic metabolism is crucial because the syntrophs can only
 506 oxidize substrates and sustain growth under anaerobic conditions if the methanogens
 507 rapidly consume the products to maintain them at very low concentrations^{51,53}. Due to
 508 the fastidious metabolism, syntrophs are usually present in low abundance, and can
 509 easily become the bottleneck in the anaerobic digestion process^{7,8}. Accordingly, we
 510 investigated the global diversity of this functional guild in the ADs sampled (Fig. 7,
 511 Supplementary Fig. 9).

512

513 A clear effect of the primary substrates and digester temperature was observed on the
 514 composition and abundance of syntrophic genera in the digesters (Fig. 7a,
 515 Supplementary Fig. 9a). The most abundant genus across substrates and temperature
 516 was *Syntrophaceticus*, despite being barely detected in ADs treating wastewater sludge.
 517 The type strain of this genus, *S. schinkii* Sp3^T, is an acetate-oxidizing syntroph that
 518 thrives, and has a competitive advantage, under high ammonium concentrations (up to
 519 8,400 mgN/L)^{54,55}. The lack of *Syntrophaceticus* in ADs treating wastewater sludge
 520 may therefore be explained by lower ammonium concentrations in these ADs
 521 (1,617±4,312 mgN/L, n=145) compared to those treating food waste (2,913±1,681
 522 mgN/L, n=33), and manure (3,449±933 mgN/L, n=18).

523



524

525 Fig. 7: Global diversity of syntrophs based on V1-V3 amplicon data. The percent relative abundance
 526 represents the mean for genera across a) different temperature range and primary substrates, and b)
 527 different countries considering only mesophilic ADs treating mainly wastewater sludge. Colored circles
 528 next to the genus labels indicate whether the genera have previously been identified as growing in ADs
 529 at Danish WWTPs according to Jiang *et al.*⁷. Blue: >50% of ASVs classified as growing; Yellow: 20-
 530 50% of ASVs classified as growing. Red: <20% of ASVs classified as growing. Gray: No information
 531 available for the specific genus.

532 *Syntrophomonas*, the second most abundant genus, was common in all AD types
533 investigated, indicating a broader ecological niche. Isolated representatives from this
534 genus can grow syntrophically via β -oxidation of saturated fatty acids of various lengths
535 (C4-C18, depending on strain) ⁵⁶⁻⁵⁹, and they are therefore likely important for the
536 conversion of long-chain fatty acids in ADs. Among the abundant syntrophs,
537 *Tepidimicrobium*, a member of the order Clostridiales, was also observed in all AD
538 types except mesophilic ADs treating wastewater sludge. The exact metabolism of
539 *Tepidimicrobium* in ADs remains to be determined, however all isolated representatives
540 can degrade proteinaceous compounds and some species can also use carbohydrates ⁶⁰.
541 Furthermore, *Tepidimicrobium* has been proposed to grow syntrophically by direct
542 interspecies electron transfer (DIET) with *Methanothermobacter* in a process like that
543 observed for *Geobacter* ⁶¹. Accordingly, it is likely that the *Tepidimicrobium* acts as a
544 syntrophic primary degrader in the ADs targeting mainly proteins, carbohydrates, and
545 derivatives.

546

547 Finally, we observed a high abundance of the genus *Smithella* in mesophilic ADs
548 treating industrial waste, manure, and wastewater sludge. The type strain *S. propionica*
549 LYP^T is a propionate oxidizing syntroph, which uses a unique dismutation pathway in
550 which propionate is first converted to acetate and butyrate, and the latter is hereafter β -
551 oxidized syntrophically to acetate and hydrogen ^{62,63}. Calculations of Gibbs free energy
552 for this special propionate metabolism indicates a higher tolerance toward elevated
553 hydrogen concentrations ⁶⁴, which could explain why some *Smithella* prevail in certain
554 ADs. However, *Smithella* has also been implicated in the syntrophic degradation of
555 long-chain alkanes ^{65,66}, which could reflect a more versatile metabolism.

556

557 When investigating geographical diversity of syntrophic fatty acid oxidizing bacteria
558 in mesophilic ADs treating wastewater sludge, a similar pattern was observed across
559 countries (Fig. 7b, Supplementary Fig. 9b). *Smithella*, was generally the dominating
560 syntroph. However, *Syntrophomonas*, *Syntrophorhabdus*, *Ca. Phosphitivorax*, and
561 *Syntrophus* also occurred at a high relative abundance in almost all countries. Isolates
562 of *Syntrophorhabdus*, including the type strain *S. aromaticus* UI^T, are syntrophic
563 fermenters of aromatic compounds and may accordingly play an important role in the
564 detoxification of these substrates in ADs ^{67,68}. *Ca. Phosphitivorax* was recently
565 discovered as a butyrate degrading syntroph by genome-resolved meta-transcriptomics
566 in a digester treating wastewater sludge ⁵², and *Syntrophus* participates in the
567 degradation of fatty acids and aromatics ^{69,70}. Overall, the results suggest a complex
568 syntrophic degradation process, which involves multiple genera with different substrate
569 specificities.

570

571 To gain additional insight into the global diversity of syntrophs, we also investigated
572 the species-level diversity across mesophilic digesters treating wastewater sludge
573 (Supplementary Fig. 10). We observed a large species diversity among most of the

574 abundant syntrophic genera. Furthermore, we found that the most abundant species in
575 the ADs were often distinct from the isolated representatives, which prompts for further
576 investigations into the metabolic potential of syntrophs *in situ*.

577

578 ***Global diversity of filamentous bacteria***

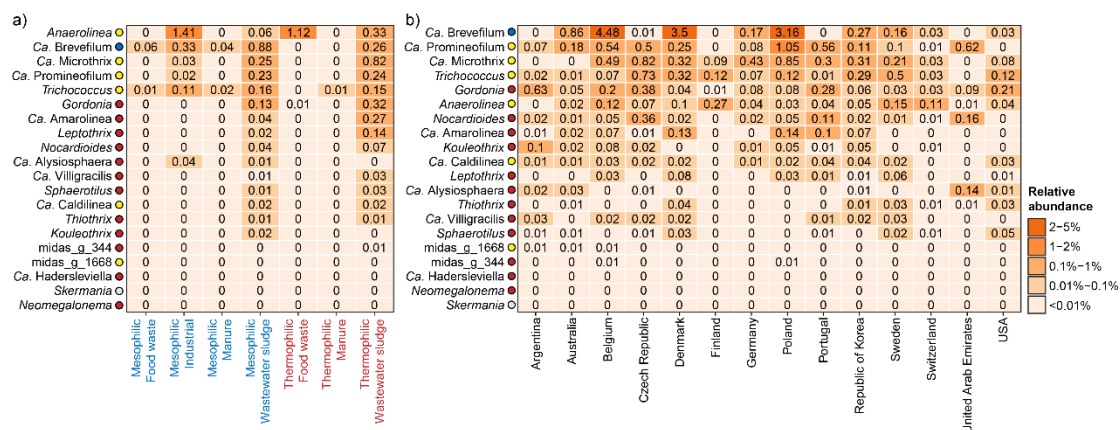
579 Foaming is a common operational problem in ADs and has a strong negative impact on
580 process performance resulting in considerable costs. Both abiotic and biotic factors are
581 involved in foaming⁷¹. The abiotic factors include high loading rates of surfactants (oil,
582 grease, fatty acids, detergent, proteins, and particulate matter) and biosurfactants
583 produced by microbes in the digester⁷². The biotic factors cover increased abundance
584 of hydrophobic, filamentous microorganisms that can interact with, and stabilize, gas
585 bubbles in the foam^{71,73}. To gain further insight into potential foam forming microbes,
586 we examined the global diversity of known filamentous bacteria in ADs (Fig. 8,
587 Supplementary Fig. 11).

588

589 The diversity and mean relative abundance of known filamentous organisms were
590 generally low in the ADs examined except for those treating wastewater sludge (Fig.
591 8a, Supplementary Fig. 11). However, the increased diversity and abundance in the
592 latter are to a large extent the result of passive immigration from the fed surplus sludge.
593 However, most of these are unable to grow in the ADs⁷. *Anaerolinea*, *Ca. Brevefilum*,
594 and *Trichococcus* were common across ADs treating all primary substrates (Fig. 8a,
595 Supplementary Fig. 11), whereas *Ca. Microthrix* and *Ca. Promineofilum* were mainly
596 observed in ADs treating wastewater sludge. Many of the Chloroflexi genera found
597 here were also observed in a recent meta-analysis of amplicon data from 17 studies
598 representing 62 ADs⁷⁴. Several of the abundant filamentous genera, including *Ca.*
599 *Microthrix* and *Ca. Brevefilum*, were previously found to correlate with the foaming
600 potential of full-scale digester sludge from mesophilic ADs at WWTPs⁷³. *Ca.*
601 *Brevefilum* seems especially interesting as it grows well in ADs^{7,75}.

602

603 The species-level diversity was generally low for the filamentous bacteria
604 (Supplementary Fig. 12). *Ca. Brevefilum* was dominated by *Ca. B. fermentans*,
605 *Trichococcus* by *midas_s_4*, *Ca. Microthrix* by *Ca. M. parvicella* and *Ca. M.*
606 *subdominans*, and *Gordonia* by *G. defluvii* and *G. amarae*. *Ca. Promineofilum* was
607 dominated by *Ca. P. glycogenico*, but a few MiDAS placeholder species, were also
608 commonly observed. The low species-level diversity of potential foam-forming
609 bacteria suggests that it may be feasible to develop and implement universal mitigation
610 strategies for these bacteria in ADs worldwide.



611

612

613

614

615

616

617

618

619

620

Fig. 8: Global diversity of known filamentous organisms based on V1-V3 amplicon data. The percent relative abundance represents the mean for genera across a) different temperature range and primary substrates, and b) different countries considering only mesophilic ADs treating mainly wastewater sludge. Colored circles next to the genus labels indicate whether the genera have previously been identified as growing in ADs at Danish WWTPs according to Jiang *et al.*⁷. Blue: >50% of ASVs classified as growing; Yellow: 20-50% of ASVs classified as growing. Red: <20% of ASVs classified as growing. Gray: No information available for the specific genus.

621

Conclusion and perspectives

622

MiDAS 5 was made possible thanks to a huge collaborative effort from experts worldwide, who contributed to the project by sampling and providing metadata for ADs in their respective countries. Building on the success of its predecessor, MiDAS 4, this latest expansion covers ASV-resolved, full-length 16S rRNA gene references from numerous ADs from all parts of the globe covering different operations parameters and different substrates. This expanded database provides greatly improved coverage for AD-specific taxa and a strongly needed taxonomy for uncultured lineages, which lacks official taxonomic classification. As such, it will be an invaluable resource for researchers and AD professionals, providing them with a common point of reference to facilitate knowledge sharing and pave the way for a comprehensive understanding of the AD microbiome.

623

624

625

626

627

628

629

630

631

632

633

Our in silico 16S rRNA gene primer evaluation based on the MiDAS 5 database revealed that the coverage of commonly applied primer pairs varies significantly, with some having low coverage and potential bias towards certain taxa. Because the primer coverage was evaluated for all taxa in the MiDAS 5 database at all taxonomic ranks, it provides a solid foundation for designing experiments and targeting specific taxa in future studies. For general microbial profiling of ADs, we would recommend the use of the newly improved universal V4 primer pair³⁸, as it show excellent coverage for both archaea and bacteria in the AD ecosystem.

634

635

636

637

638

639

640

641 Although the total microbial diversity in ADs is huge, importantly, we showed that less
642 than 1000 genera and species accounted for most of the microbes in the AD ecosystem.
643 By focusing on the fraction of these abundant and common microbes that can grow in
644 the AD systems, we will be able to explain most of the microbial processes that occur
645 in the anaerobic digestion process. This list of “Most Wanted” organisms contain
646 species that should be prime targets for future in situ studies and the reconstruction of
647 MAGs. These genomes can then be annotated to provide additional details about their
648 potential metabolic pathways and roles in the AD ecosystem^{15,16,76–78}.

649

650 The global survey of the ADs microbiota using three different primer pairs provided a
651 unique insight into the global diversity of individual AD taxa and clues into the
652 environmental and operational factors that define their ecological niches. This
653 information will be invaluable in the development of future microbiome management
654 strategies and improved sustainability of the field of anaerobic digestion.

655

656 To enhance knowledge dissemination, we have updated the MiDAS Field Guide
657 available at www.midasfieldguide.org. This dynamic resource allows users to delve
658 into specifics related to the physiology, morphology, and ecology of genera listed in
659 the MiDAS database. Additionally, it offers country-specific data on the prevalence of
660 all MiDAS genera and species in WWTPs and ADs. Finally, it provides information on
661 the availability of fluorescence in situ hybridization probes and reference genomes,
662 paving the way for subsequent research endeavors.

663 **Methods:**

664 ***Sampling and metadata collection***

665 To facilitate sampling of ADs worldwide, we established the MiDAS Global
666 Consortium for Anaerobic Digesters, which consists of 25 anaerobic digestion experts
667 in 19 countries. Members of the consortium acted as national sampling coordinators
668 and were in direct contact with the ADs. Two samples were obtained from each AD
669 and shipped on ice to the sampling coordinators. For each replicate, 2 mL sample was
670 preserved in 2 mL RNAlater (Invitrogen), stored at 4°C until all national samples were
671 collected (usually within a few days), and then shipped to Aalborg University with
672 cooling elements. Upon arrival, the samples were separated into aliquots that were
673 prepared for nucleic acid purification. Metadata associated with each AD was also
674 obtained by the sampling coordinators and is provided as Supplementary Data 1.
675 Minimum information from all ADs included continent, country, GPS coordinates,
676 sampling date, temperature in the digester ("Mesophilic" ($\leq 45^\circ\text{C}$) or "Thermophilic"
677 ($50\text{--}60^\circ\text{C}$)), primary substrate ("Wastewater sludge", "Industrial", "Food waste",
678 "Manure", or "Other"), and digester technology ("Two-stage digester (TSAD)",
679 "Continuous Stirred Tank Reactors (CSTR)", "Upflow anaerobic sludge blanket
680 (UASB)", or "Other").

681

682 ***General molecular methods***

683 All commercial kits were used according to the protocols provided by the manufacturer
684 unless otherwise stated. The concentration and quality of nucleic acids were determined
685 using a Qubit 3.0 fluorometer (Thermo Fisher Scientific) and an Agilent 2200
686 TapeStation (Agilent Technologies), respectively.

687

688 ***Nucleic acid purification***

689 DNA was purified using a custom plate-based extraction protocol based on the
690 FastDNA spin kit for soil (MP Biomedicals). The protocol is available at
691 www.midasfieldguide.org (aau_ad_dna_v 2.0). RNAlater preserved samples were
692 thawed and homogenized using a Heidolph RZR 2020 laboratory stirrer. 20 μL of
693 sample was resuspended in 300 μL PBS and transferred to Lysing Matrix E barcoded
694 tubes (MP Biomedicals). 40 μL of MT buffer was added and lysis was performed by
695 bead beating in a FastPrep-96 bead beater (MP Biomedicals) (3x 120 s, 1800 rpm with
696 2 min incubation on ice between cycles). The samples were centrifuged ($3,486 \times g$, 10
697 min) and 200 μL supernatant was transferred to a 96-well PCR-plate. 50 μL Protein
698 Precipitation Solution (PPS) was mixed with each sample, which was then centrifuged
699 again. 150 μL supernatant was cleaned-up using 100 μL CleanNGS beads with elution
700 into 60 μL of nuclease-free water. 40 μL of the purified DNA was transferred to a new
701 96-well plate and stored at -80°C .

702 ***Full-length 16S rRNA gene library preparation, sequencing, and processing***

703 Full-length 16S rRNA gene sequencing was carried out using high-accuracy, long-read
704 amplicon sequencing using unique molecular identifiers (UMIs) and PacBio circular
705 consensus sequencing (CCS)⁷⁹. Oligonucleotides used can be found in Supplementary
706 Table 1. Bacterial and archaeal 16S rRNA genes were UMI-tagged using overhang
707 primers based on the 27F and 1391R⁸⁰ and SSU1ArF and SSU1000ArR³⁴ primer pairs,
708 respectively. These primers have shown excellent coverage for the known bacterial and
709 archaeal diversity in silico^{34,80}.

710

711 *Addition of UMI-tags by overhang PCR:* Adaptors containing UMIs, and defined
712 primer binding sites were added to each end of the bacterial and archaeal 16S rRNA
713 genes by PCR. The reaction contained 20 µL of 5x SuperFi Buffer (Invitrogen), 2 µL
714 of 10 mM dNTP mix, 5 µL of 10 µM f16S_pcr1_fw, 5 µL of 10 µM f16S_pcr1_rv, 1
715 µL of 2 U/µL Platinum SuperFi DNA polymerase (Invitrogen), 100 ng of pooled
716 template DNA (from all ADs), and nuclease-free water to 100 µL. The reaction was
717 incubated with an initial denaturation at 98°C for 30 s followed by 2 cycles of
718 denaturation at 98°C for 20 s, annealing at 55°C for 30 s, and extension at 72°C for 45
719 s, and then a final extension at 72°C for 5 min. The sample was purified using 0.6x
720 CleanNGS beads and eluted in 20 µL nuclease-free water.

721 *Primary library amplification:* The tagged 16S rRNA gene amplicons were amplified
722 using PCR to obtain enough product for quantification. The reaction contained 19 µL
723 of UMI-tagged sample, 20 µL 5x SuperFi buffer (Invitrogen), 2 µL of 10 mM dNTP, 5
724 µL of 10 µM f16S_pcr2_fw, 5 µL of 10 µM f16S_pcr2_rv, 48 µL nuclease-free water,
725 and 1 µL 2U/µL Platinum SuperFi DNA polymerase (Invitrogen). The reaction was
726 incubated with an initial denaturation at 98°C for 30 s followed by 15 cycles of
727 denaturation at 98°C for 20 s, annealing at 60°C for 30 s, and extension at 72°C for 45
728 s and then a final extension at 72°C for 5 min. The PCR product was purified using
729 0.6x CleanNGS beads and eluted in 11 µL nuclease-free water. The amplicons were
730 validated on a Genomic screentape and quantified with the Qubit dsDNA HS assay kit.

731 *Clonal library amplification:* Tagged amplicon libraries were diluted to approximately
732 250,000 molecules/µL and amplified by PCR to obtain clonal copies of each uniquely
733 tagged amplicon molecule. Three libraries were made for the bacterial 16S rRNA genes
734 and one for archaea. The PCR reactions contained 1 µL diluted primary library, 20 µL
735 5x SuperFi buffer (Invitrogen), 2 µL of 10 mM dNTP, 5 µL of 10 µM f16S_pcr2_fw,
736 5 µL of 10 µM f16S_pcr2_rv, 66 µL nuclease-free water, and 1 µL 2U/µL Platinum
737 SuperFi DNA polymerase (Invitrogen). The reaction was incubated with an initial
738 denaturation at 98°C for 30 s followed by 25 cycles of denaturation at 98°C for 20 s,
739 annealing at 60°C for 30 s, and extension at 72°C for 45 s and then a final extension at
740 72°C for 5 min. The PCR product was purified using 0.6x CleanNGS beads and eluted
741 in 20 µL nuclease-free water. The amplicons were validated on a Genomic screentape
742 and quantified with the Qubit dsDNA HS assay kit.

743 *PacBio CCS sequencing:* The four clonal libraries were sent to Admera Health
744 (Plainfield, NJ, USA) for PacBio library preparation and sequencing. Here amplicons
745 were incubated with T4 polynucleotide kinase (New England Biolabs) following the
746 manufacturer's instructions, and sequencing library prepared using SMRTbell Express
747 Template Preparation kit 1.0 following the standard protocol. Sequencing was
748 performed using 4x SMRT cells on a Sequel II using a Sequel II Sequencing kit 1.0,
749 Sequel II Binding and Int Ctrl kit 1.0 and Sequel II SMRT Cell 8M, following the
750 standard protocol with 1 h pre-extension and 15 h collection time (Pacific Biosciences).

751 *Bioinformatic processing:* CCS reads were generated from raw PacBio data using CCS
752 v.3.4.1 (<https://github.com/PacificBiosciences/ccs>) with default settings. UMI
753 consensus sequences (consensus_raconx3.fa) were obtained using the longread_umi
754 script (https://github.com/SorenKarst/longread_umi)⁷⁹ using the following options:
755 pacbio_pipeline, -v 3, -m 1000, -M 2000, -s 60, -e 60, -f
756 CAAGCAGAAGACGGCATAACGAGAT, -F AGRGTTYGATYMTGGCTCAG
757 (bacteria) or TCCGGTTGATCCYGCBRG (archaea), -r
758 AATGATACGGCGACCACCGAGATC, -R GACGGGCGGTGWGTRCA (bacteria)
759 or GGCCATGCAMYWCCTCTC (archaea), and -c 3. The UMI-consensus reads were
760 oriented based on the SILVA 138.1 SSURef NR99 database using the usearch
761 v.11.0.667 -orient command and trimmed between the 27f and 1391r (bacteria) or
762 SSU1ArF and SSU1000ArR (archaea) primer binding sites using the trimming function
763 in CLC genomics workbench v. 20.0. Sequences without both primer binding sites were
764 discarded. The trimmed high-fidelity reads were processed with AutoTax v. 1.7.4⁶ to
765 create full-length 16S rRNA gene amplicon sequence variants (FL-ASV) and these
766 were added to the MiDAS 4.8.1 reference database²⁰ to create MiDAS 5.0. Subsequent
767 updates to MiDAS 5.2 were made to accommodate taxonomic updates (see the release
768 change logs for details).

769 *Short-read amplicon sequencing*

770 V1-V3 amplicons were made using the 27F (5'-AGAGTTTGATCCTGGCTCAG-3')
771 ⁸¹ and 534R (5'-ATTACCGCGGCTGCTGG-3')⁸² primers with barcodes and Illumina
772 adaptors (IDT)⁸³. 25 µL PCR reactions in duplicate were run for each sample using 1X
773 PCR BIO Ultra Mix (PCR Biosystems), 400 nM of both forward and reverse primer,
774 and 10 ng template DNA. PCR conditions were 95°C, for 2 min followed by 20 cycles
775 of 95°C for 20 s, 56°C for 30 s, and 72°C for 60 s, followed by a final elongation at
776 72°C for 5 min. PCR products were purified using 0.8x CleanNGS beads and eluted in
777 25 µL nuclease-free water.

778
779 V3-V5 amplicons were made using the Arch-340F (5'-CCCTAHGGGGYGCASCA-
780 3') and Arch-915R (5'-GWGCYCCCCCGYCAATTC-3') primers⁸⁴. 25 µL PCR
781 reactions in duplicate were run for each sample using 1X PCR BIO Ultra Mix (PCR
782 Biosystems), 400 nM of both forward and reverse primer, and 10 ng template DNA.
783 PCR conditions were 95°C, for 2 min followed by 30 cycles of 95°C for 15 s, 55°C for

784 15 s, and 72°C for 50 s, followed by a final elongation at 72°C for 5 min. PCR products
785 were purified using 0.8x CleanNGS beads and eluted in 25 µL nuclease-free water. 2
786 µL of purified PCR product from above was used as template for a 25 µL Illumina
787 barcoding PCR reaction containing 1x PCR BIO Reaction buffer, 1 U PCR BIO HiFi
788 Polymerase (PCR Biosystems) and 10 µL of Nextera adaptor mix (Illumina). PCR
789 conditions were 95°C, for 2 min, 8 cycles of 95°C for 20 s, 55°C for 30 s, and 72°C for
790 60 s, followed by a final elongation at 72°C for 5 min. PCR products were purified
791 using 0.8x CleanNGS beads and eluted in 25 µL nuclease-free water.

792

793 V4 amplicons were made using the 515F (5'-GTGYCAGCMGCCGCGGTAA-3')⁸²
794 and 806R (5'-GGACTACNVGGGTWTCTAAT-3')⁸⁵ primers. 25 µL PCR reactions
795 in duplicate were run for each sample using 1X PCR BIO Ultra Mix (PCR Biosystems),
796 400 nM of both forward and reverse primer, and 10 ng template DNA. PCR conditions
797 were 95°C, for 2 min followed by 30 cycles of 95°C for 15 s, 55°C for 15 s, and 72°C
798 for 50 s, followed by a final elongation at 72°C for 5 min. PCR products were purified
799 using 0.8x CleanNGS beads and eluted in 25 µL nuclease-free water. 2 µL of purified
800 PCR product from above was used as template for a 25 µL Illumina barcoding PCR
801 reaction as described for the V3-V5 amplicons.

802

803 16S rRNA gene V1-V3, V3-V5, and V4 amplicon libraries were pooled separately in
804 equimolar concentrations and diluted to 4 nM. The amplicon libraries were paired-end
805 sequenced (2 × 300 bp) on the Illumina MiSeq using v3 chemistry (Illumina, USA). 10
806 to 20% PhiX control library was added to mitigate low diversity library effects.

807

808 ***Processing of short-read amplicon data***

809 Usearch v.11.0.667⁸⁶ was used for processing of 16S rRNA gene amplicon data and
810 for read mapping. V1-V3 forward and reverse reads were merged using the `usearch -`
811 `fastq_mergepairs` command, filtered to remove phiX sequences using `usearch -`
812 `filter_phix`, and quality filtered using `usearch -fastq_filter` with `-fastq_maxee 1.0`.
813 Dereplication was performed using `-fastx_uniques` with `-sizeout`, and amplicon
814 sequence variants (ASVs) were resolved using the `usearch -unnoise3` command⁸⁷. An
815 ASV-table was created by mapping the quality filtered reads to the ASVs using the
816 `usearch -otutab` command with the `-zotus` and `-strand plus` options. Taxonomy was
817 assigned to ASVs using the `usearch -sintax` command with `-strand both` and `-`
818 `sintax_cutoff 0.8` options. Mapping of ASVs to reference databases was done with the
819 `usearch -usearch_global` command and the `-id 0`, `-maxaccepts 0`, `-maxrejects 0`, `-`
820 `top_hit_only`, and `-strand plus` options.

821

822 16S rRNA gene V3-V5 forward reads (reverse reads in relation the 16S rRNA gene)
823 were filtered to remove phiX sequences using `usearch -filter_phix`, trimmed to remove
824 primers and obtain a fixed length of 250 bp using `-fastx_truncate` with `-stripleft -17` and
825 `truncLen 250`, reverse complemented with `usearch -fastx_revcomp`, and quality filtered

826 using usearch -fastq_filter with -fastq_maxee 1.0. Subsequent processing was like that
827 for the V1-V3 amplicons.

828

829 16S rRNA gene V4 forward reads (reverse reads in relation the 16S rRNA gene) were
830 trimmed with cutadapt v.2.8⁸⁸ based on the V4 primers with the -g
831 ^GGACTACHVGGGTWTCTAAT...TTACCGCGGCKGCTGGCAC and --discard-
832 untrimmed options. The trimmed reads, which span the entire V4 amplicon, were
833 reverse complemented with usearch -fastx_revcomp, and quality filtered using usearch
834 -fastq_filter with -fastq_maxee 1.0. Subsequent processing was like that for the V1-V3
835 amplicons.

836

837 ***In silico primer evaluation***

838 The specificity of commonly used amplicon primers was determined for each FL-ASV
839 using the analyze_primers.py script from Primer Prospector v. 1.0.1⁸⁹. The specificity
840 of primer sets was defined based on the overall weighted scores (OWS) for the primer
841 with the highest score as follows: Perfect hit (OWS = 0), partial hit (OWS > 0, and ≤
842 1), poor hit (OWS > 1). The percentage of perfect hits were calculated in R for all taxa
843 in MiDAS 5.

844

845 ***Microbial community analyses***

846 Short-read amplicon data was analyzed with R v.4.3.2⁹⁰ through RStudio IDE
847 v.2023.12.1⁹¹, with the tidyverse v.2.0.0 (<https://www.tidyverse.org/>), vegan v.2.6-4
848⁹², maps v.3.4.2⁹³, data.table v.1.14.10⁹⁴, FSA v.0.9.5⁹⁵, rcompanion v. 2.4.35⁹⁶,
849 patchwork v.1.1.3⁹⁷, ggupset v.0.3.0⁹⁸ and Ampvis2 v.2.8.6⁹⁹ packages.

850

851 The microbial community analyses were performed based on all three 16S rRNA gene
852 short-read amplicon dataset (V1-V3, V3-V5, and V4). Samples with less than 10,000
853 reads and those lacking information about digester technology, primary substrate, and
854 temperature in the digester were discarded from the analyses. After filtration, 547 V1-
855 V3, 542 V3-V5, and 430 V4 samples remained.

856 Associations between the AD microbiota and the following process-related or
857 environmental variables were investigated: Digester technology, primary substrate,
858 temperature in the digester, and continent (see definitions above). All variables were
859 treated as factors.

860 For alpha diversity analyses, samples were rarefied to 10,000 reads, and alpha diversity
861 (observed ASVs and inverse Simpson) was calculated using the ampvis2 package. The
862 Kruskal-Wallis with Dunn's post-hoc test (Bonferroni correction with $\alpha=0.01$ before
863 correction) was used to determine statistically significant differences in alpha diversity
864 between samples grouped by process and environmental variables.

865 Beta diversity distances based on Bray-Curtis (abundance-based) for genera was
866 calculated using the `vegdist` function in the `vegan` R package and visualized by PCoA
867 plots with the `ampvis2` package. To determine how much individual parameters
868 affected the structure of the microbial community across the ADs, a permutational
869 multivariate analysis of variance (PERMANOVA) test was performed on the beta-
870 diversity matrices using the `adonis` function in the `vegan` package with 999
871 permutations.

872 Core taxa (genera and species) were determined separately for ADs treating different
873 primary substrates and operating at different temperatures (mesophilic and
874 thermophilic) based on their relative abundances in individual ADs according to the
875 three short-read amplicon datasets. Core taxa definitions were identical to those applied
876 in the MiDAS global survey of WWTPs²⁰. Taxa were considered abundant when
877 present at >0.1% relative read abundance in individual ADs. Based on how frequently
878 taxa were observed to be abundant, we defined the following core communities: loose
879 core (>20% of ADs), general core (>50% of ADs), and strict core (>80% of ADs).
880 Additionally, we defined conditionally rare or abundant taxa (CRAT)¹⁰⁰ composed of
881 taxa present in one or more ADs at >1% relative abundance, but not belonging to the
882 core taxa.

883

884 **Data Availability:**

885 The raw and assembled sequencing data generated in this study have been deposited in
886 the NCBI SRA database under accession code PRJNA1019951
887 [<https://www.ncbi.nlm.nih.gov/bioproject/PRJNA1019951>]. The MiDAS 5 reference
888 database in SINTAX, QIIME and DADA2 format is available at the MiDAS fieldguide
889 website [<https://www.midasfieldguide.org/guide/downloads>].

890

891 **Code Availability:**

892 R scripts used for data analyses and figures are available at GitHub
893 [<https://github.com/msduholm/MiDAS5>]¹⁰¹. Raw data files for the R scripts are
894 available at Figshare [<https://doi.org/10.6084/m9.figshare.24219199.v2>]¹⁰².

895 **References:**

- 896 1. Tiwary, A., Williams, I. D., Pant, D. C. & Kishore, V. V. N. Emerging perspectives
897 on environmental burden minimisation initiatives from anaerobic digestion
898 technologies for community scale biomass valorisation. *Renewable and*
899 *Sustainable Energy Reviews* **42**, 883–901 (2015).
- 900 2. Achinas, S., Achinas, V., & Gerrit Jan Willem Euverink. A technological overview
901 of biogas production from biowaste. *Engineering* **3**, 299–307 (2017).
- 902 3. Samoraj, M. *et al.* The challenges and perspectives for anaerobic digestion of
903 animal waste and fertilizer application of the digestate. *Chemosphere* **295**, 133799
904 (2022).
- 905 4. Czekala, W., Jasiński, T., Grzelak, M., Witaszek, K. & Dach, J. Biogas plant
906 operation: Digestate as the valuable product. *Energies* **15**, 8275 (2022).
- 907 5. Briones, A. & Raskin, L. Diversity and dynamics of microbial communities in
908 engineered environments and their implications for process stability. *Current*
909 *Opinion in Biotechnology* **14**, 270–276 (2003).
- 910 6. Dueholm, M. S. *et al.* Generation of comprehensive ecosystem-specific reference
911 databases with species-level resolution by high-throughput full-length 16S rRNA
912 gene sequencing and automated taxonomy assignment (AutoTax). *mBio* **11**,
913 e01557-20 (2020).
- 914 7. Jiang, C. *et al.* Characterizing the growing microorganisms at species level in 46
915 anaerobic digesters at Danish wastewater treatment plants: A six-year survey on
916 microbial community structure and key drivers. *Water Research* **193**, 116871
917 (2021).
- 918 8. Vanwonterghem, I. *et al.* Deterministic processes guide long-term synchronised
919 population dynamics in replicate anaerobic digesters. *ISME Journal* **8**, 2015–2028
920 (2014).
- 921 9. Ofiteiru, I. D. *et al.* Combined niche and neutral effects in a microbial wastewater
922 treatment community. *Proceedings of the National Academy of Sciences* **107**,
923 15345–15350 (2010).
- 924 10. Kirkegaard, R. H. *et al.* The impact of immigration on microbial community
925 composition in full-scale anaerobic digesters. *Scientific Reports* **7**, 9343 (2017).
- 926 11. Werner, J. J. *et al.* Bacterial community structures are unique and resilient in full-
927 scale bioenergy systems. *Proceedings of the National Academy of Sciences of the*
928 *United States of America* **108**, 4158–4163 (2011).
- 929 12. Campanaro, S. *et al.* Metagenomic analysis and functional characterization of the
930 biogas microbiome using high throughput shotgun sequencing and a novel binning
931 strategy. *Biotechnology for biofuels* **9**, 26 (2016).
- 932 13. Calusinska, M. *et al.* A year of monitoring 20 mesophilic full-scale bioreactors
933 reveals the existence of stable but different core microbiomes in bio-waste and
934 wastewater anaerobic digestion systems. *Biotechnology for Biofuels* **11**, 196
935 (2018).

- 936 14. Mei, R. *et al.* Operation-driven heterogeneity and overlooked feed-associated
937 populations in global anaerobic digester microbiome. *Water Research* **124**, 77–84
938 (2017).
- 939 15. Ma, S. *et al.* A microbial gene catalog of anaerobic digestion from full-scale biogas
940 plants. *GigaScience* **10**, giaa164 (2021).
- 941 16. Campanaro, S. *et al.* New insights from the biogas microbiome by comprehensive
942 genome-resolved metagenomics of nearly 1600 species originating from multiple
943 anaerobic digesters. *Biotechnology for Biofuels* **13**, 25 (2020).
- 944 17. McIlroy, S. J. *et al.* MiDAS 2.0: An ecosystem-specific taxonomy and online
945 database for the organisms of wastewater treatment systems expanded for anaerobic
946 digester groups. *Database* **2017**, bax016 (2017).
- 947 18. Nierychlo, M. *et al.* MiDAS 3: An ecosystem-specific reference database,
948 taxonomy and knowledge platform for activated sludge and anaerobic digesters
949 reveals species-level microbiome composition of activated sludge. *Water Research*
950 **182**, 115955 (2020).
- 951 19. McIlroy, S. J. *et al.* MiDAS: The field guide to the microbes of activated sludge.
952 *Database : the journal of biological databases and curation* **2015**, bav062 (2015).
- 953 20. Dueholm, M. K. D. *et al.* MiDAS 4: A global catalogue of full-length 16S rRNA
954 gene sequences and taxonomy for studies of bacterial communities in wastewater
955 treatment plants. *Nat Commun* **13**, 1908 (2022).
- 956 21. Albertsen, M., Karst, S. M., Ziegler, A. S., Kirkegaard, R. H. & Nielsen, P. H. Back
957 to basics - the influence of DNA extraction and primer choice on phylogenetic
958 analysis of activated sludge communities. *PloS one* **10**, e0132783 (2015).
- 959 22. Kristensen, J. M., Singleton, C., Clegg, L.-A., Petriglieri, F. & Nielsen, P. H. High
960 diversity and functional potential of undescribed “Acidobacteriota” in Danish
961 wastewater treatment plants. *Front. Microbiol.* **12**, 906 (2021).
- 962 23. Nierychlo, M. *et al.* Low global diversity of *Candidatus* Microthrix, a troublesome
963 filamentous organism in full-scale WWTPs. *Front Microbiol* **12**, 690251 (2021).
- 964 24. Petriglieri, F. *et al.* “*Candidatus* Dechloromonas phosphoritropha” and “*Ca.* D.
965 phosphorivorans”, novel polyphosphate accumulating organisms abundant in
966 wastewater treatment systems. *ISME J* **15**, 3605–3614 (2021).
- 967 25. Petriglieri, F. *et al.* Reevaluation of the phylogenetic diversity and global
968 distribution of the genus “*Candidatus* Accumulibacter”. *mSystems* **7**, e00016-22
969 (2022).
- 970 26. Singleton, C. M. *et al.* The novel genus, ‘*Candidatus* Phosphoribacter’, previously
971 identified as *Tetrasphaera*, is the dominant polyphosphate accumulating lineage in
972 EBPR wastewater treatment plants worldwide. *ISME J* **16**, 1605–1616 (2022).
- 973 27. Yarza, P. *et al.* Uniting the classification of cultured and uncultured bacteria and
974 archaea using 16S rRNA gene sequences. *Nature Reviews Microbiology* **12**, 635–
975 645 (2014).
- 976 28. Rivière, D. *et al.* Towards the definition of a core of microorganisms involved in
977 anaerobic digestion of sludge. *ISME Journal* **3**, 700–714 (2009).

- 978 29. Fujimoto, M., Carey, D. E., Zitomer, D. H. & McNamara, P. J. Syntroph diversity
979 and abundance in anaerobic digestion revealed through a comparative core
980 microbiome approach. *Appl Microbiol Biotechnol* **103**, 6353–6367 (2019).
- 981 30. Westerholm, M., Calusinska, M. & Dolfing, J. Syntrophic propionate-oxidizing
982 bacteria in methanogenic systems. *FEMS Microbiology Reviews* **46**, fuab057
983 (2022).
- 984 31. Johnson, J. S. *et al.* Evaluation of 16S rRNA gene sequencing for species and strain-
985 level microbiome analysis. *Nature Communications* **10**, 1–11 (2019).
- 986 32. McDonald, D. *et al.* Greengenes2 unifies microbial data in a single reference tree.
987 *Nat Biotechnol* 1–4 (2023) doi:10.1038/s41587-023-01845-1.
- 988 33. Gonzalez, A. *et al.* Qiita: rapid, web-enabled microbiome meta-analysis. *Nat*
989 *Methods* **15**, 796–798 (2018).
- 990 34. Bahram, M., Anslan, S., Hildebrand, F., Bork, P. & Tedersoo, L. Newly designed
991 16S rRNA metabarcoding primers amplify diverse and novel archaeal taxa from
992 the environment. *Environmental Microbiology Reports* **11**, 487–494 (2018).
- 993 35. Pausan, M. R. *et al.* Exploring the archaeome: Detection of archaeal signatures in
994 the human body. *Frontiers in Microbiology* **10**, 2796 (2019).
- 995 36. Wu, L. *et al.* Global diversity and biogeography of bacterial communities in
996 wastewater treatment plants. *Nature Microbiology* **4**, 1183–1195 (2019).
- 997 37. Thompson, L. R. *et al.* A communal catalogue reveals Earth’s multiscale microbial
998 diversity. *Nature* **551**, 457–463 (2017).
- 999 38. Hu, H. *et al.* Global abundance patterns, diversity, and ecology of Patescibacteria
1000 in wastewater treatment plants. *Microbiome* **12**, 55 (2024).
- 1001 39. Mei, R., Narihito, T., Nobu, M. K., Kuroda, K. & Liu, W.-T. Evaluating digestion
1002 efficiency in full-scale anaerobic digesters by identifying active microbial
1003 populations through the lens of microbial activity. *Scientific Reports* **6**, 34090
1004 (2016).
- 1005 40. Sun, L., Pope, P. B., Eijsink, V. G. H. & Schnürer, A. Characterization of microbial
1006 community structure during continuous anaerobic digestion of straw and cow
1007 manure. *Microbial Biotechnology* **8**, 815–827 (2015).
- 1008 41. Karakashev, D., Batstone, D. J. & Angelidaki, I. Influence of environmental
1009 conditions on methanogenic compositions in anaerobic biogas reactors. *Applied*
1010 *and Environmental Microbiology* **71**, 331–338 (2005).
- 1011 42. Levén, L., Eriksson, A. R. B. & Schnürer, A. Effect of process temperature on
1012 bacterial and archaeal communities in two methanogenic bioreactors treating
1013 organic household waste. *FEMS Microbiology Ecology* **59**, 683–693 (2007).
- 1014 43. Martiny, J. B. H., Jones, S. E., Lennon, J. T. & Martiny, A. C. Microbiomes in light
1015 of traits: A phylogenetic perspective. *Science* **350**, aac9323 (2015).
- 1016 44. Knights, D. *et al.* Rethinking “enterotypes”. *Cell Host & Microbe* **16**, 433–437
1017 (2014).

- 1018 45. Saunders, A. M., Albertsen, M., Vollertsen, J. & Nielsen, P. H. The activated sludge
1019 ecosystem contains a core community of abundant organisms. *The ISME journal*
1020 **10**, 11–20 (2016).
- 1021 46. Schorn, S. *et al.* Diverse methylotrophic methanogenic archaea cause high methane
1022 emissions from seagrass meadows. *Proceedings of the National Academy of*
1023 *Sciences* **119**, e2106628119 (2022).
- 1024 47. Evans, P. N. *et al.* Methane metabolism in the archaeal phylum Bathyarchaeota
1025 revealed by genome-centric metagenomics. *Science* **350**, 434–438 (2015).
- 1026 48. Nobu, M. K., Narihiro, T., Kuroda, K., Mei, R. & Liu, W. T. Chasing the elusive
1027 Euryarchaeota class WSA2: genomes reveal a uniquely fastidious methyl-reducing
1028 methanogen. *The ISME Journal* **10**, 2478–2487 (2016).
- 1029 49. Jung, M.-Y. *et al.* A hydrophobic ammonia-oxidizing archaeon of the
1030 *Nitrosocosmicus* clade isolated from coal tar-contaminated sediment.
1031 *Environmental Microbiology Reports* **8**, 983–992 (2016).
- 1032 50. Huang, W.-C. *et al.* Comparative genomic analysis reveals metabolic flexibility of
1033 Woesearchaeota. *Nat Commun* **12**, 5281 (2021).
- 1034 51. Morris, B. E. L., Henneberger, R., Huber, H. & Moissl-Eichinger, C. Microbial
1035 syntrophy: Interaction for the common good. *FEMS Microbiology Reviews* **37**,
1036 384–406 (2013).
- 1037 52. Hao, L. *et al.* Novel syntrophic bacteria in full-scale anaerobic digesters revealed
1038 by genome-centric metatranscriptomics. *ISME J* **14**, 906–918 (2020).
- 1039 53. Sieber, J. R., McInerney, M. J. & Gunsalus, R. P. Genomic insights into syntrophy:
1040 the paradigm for anaerobic metabolic cooperation. *Annual Review of Microbiology*
1041 **66**, 429–452 (2012).
- 1042 54. Westerholm, M., Roos, S. & Schnürer, A. *Syntrophaceticus schinkiigen* nov., sp.
1043 nov., an anaerobic, syntrophic acetate-oxidizing bacterium isolated from a
1044 mesophilic anaerobic filter. *FEMS Microbiology Letters* **309**, 100–104 (2010).
- 1045 55. Schnürer, A., Müller, B. & Westerholm, M. *Syntrophaceticus*. in *Bergey's Manual*
1046 *of Systematics of Archaea and Bacteria* 1–10 (John Wiley & Sons, Ltd, 2018).
1047 doi:10.1002/9781118960608.gbm01452.
- 1048 56. McInerney, M. J., Bryant, M. P., Hespell, R. B. & Costerton, J. W. *Syntrophomonas*
1049 *wolfei* gen. nov. sp. nov., an anaerobic, syntrophic, fatty acid-oxidizing bacterium.
1050 *Applied and Environmental Microbiology* **41**, 1029–1039 (1981).
- 1051 57. Wu, C., Liu, X. & Dong, X. *Syntrophomonas cellicola* sp. nov., a spore-forming
1052 syntrophic bacterium isolated from a distilled-spirit-fermenting cellar, and
1053 assignment of *Syntrophospora bryantii* to *Syntrophomonas bryantii* comb. nov. *Int*
1054 *J Syst Evol Microbiol* **56**, 2331–2335 (2006).
- 1055 58. Sousa, D. Z., Smidt, H., Alves, M. M. & Stams, A. J. M. *Syntrophomonas zehnderi*
1056 sp. nov., an anaerobe that degrades long-chain fatty acids in co-culture with
1057 *Methanobacterium formicicum*. *Int J Syst Evol Microbiol* **57**, 609–615 (2007).
- 1058 59. Hatamoto, M., Imachi, H., Fukayo, S., Ohashi, A. & Harada, H. *Syntrophomonas*
1059 *palmitatica* sp. nov., an anaerobic, syntrophic, long-chain fatty-acid-oxidizing

- 1060 bacterium isolated from methanogenic sludge. *International Journal of Systematic*
1061 *and Evolutionary Microbiology* **57**, 2137–2142 (2007).
- 1062 60. Niu, L., Song, L., Liu, X. & Dong, X. *Tepidimicrobium xylanilyticum* sp. nov., an
1063 anaerobic xylanolytic bacterium, and emended description of the genus
1064 *Tepidimicrobium*. *International Journal of Systematic and Evolutionary*
1065 *Microbiology* **59**, 2698–2701 (2009).
- 1066 61. Wang, G., Li, Q., Gao, X. & Wang, X. C. Sawdust-derived biochar much mitigates
1067 VFAs accumulation and improves microbial activities to enhance methane
1068 production in thermophilic anaerobic digestion. *ACS Sustainable Chem. Eng.* **7**,
1069 2141–2150 (2019).
- 1070 62. Liu, Y., Balkwill, D. L., Aldrich, H. C., Drake, G. R. & Boone, D. R.
1071 Characterization of the anaerobic propionate-degrading syntrophs *Smithella*
1072 *propionica* gen. nov., sp. nov. and *Syntrophobacter wolinii*. *International Journal*
1073 *of Systematic and Evolutionary Microbiology* **49**, 545–556 (1999).
- 1074 63. de Bok, F. A. M., Stams, A. J. M., Dijkema, C. & Boone, D. R. Pathway of
1075 propionate oxidation by a syntrophic culture of *Smithella propionica* and
1076 *Methanospirillum hungatei*. *Applied and Environmental Microbiology* **67**, 1800–
1077 1804 (2001).
- 1078 64. Dolfing, J. Syntrophic propionate oxidation via butyrate: a novel window of
1079 opportunity under methanogenic conditions. *Applied and Environmental*
1080 *Microbiology* **79**, 4515–4516 (2013).
- 1081 65. Embree, M., Nagarajan, H., Movahedi, N., Chitsaz, H. & Zengler, K. Single-cell
1082 genome and metatranscriptome sequencing reveal metabolic interactions of an
1083 alkane-degrading methanogenic community. *ISME J* **8**, 757–767 (2014).
- 1084 66. Tan, B., Nesbø, C. & Foght, J. Re-analysis of omics data indicates *Smithella* may
1085 degrade alkanes by addition to fumarate under methanogenic conditions. *ISME*
1086 *Journal* **8**, 2353–2356 (2014).
- 1087 67. Nobu, M. K. *et al.* The genome of *Syntrophorhabdus aromaticivorans* strain UI
1088 provides new insights for syntrophic aromatic compound metabolism and electron
1089 flow. *Environmental Microbiology* **17**, 4861–4872 (2015).
- 1090 68. Qiu, Y.-L. *et al.* *Syntrophorhabdus aromaticivorans* gen. nov., sp. nov., the first
1091 cultured anaerobe capable of degrading phenol to acetate in obligate syntrophic
1092 associations with a hydrogenotrophic methanogen. *Appl Environ Microbiol* **74**,
1093 2051–2058 (2008).
- 1094 69. McInerney, M. J. *et al.* The genome of *Syntrophus aciditrophicus*: Life at the
1095 thermodynamic limit of microbial growth. *Proceedings of the National Academy of*
1096 *Sciences* **104**, 7600–7605 (2007).
- 1097 70. Jackson, B. E., Bhupathiraju, V. K., Tanner, R. S., Woese, C. R. & McInerney, M.
1098 J. *Syntrophus aciditrophicus* sp. nov., a new anaerobic bacterium that degrades fatty
1099 acids and benzoate in syntrophic association with hydrogen-using microorganisms.
1100 *Arch Microbiol* **171**, 107–114 (1999).

- 1101 71. Ganidi, N., Tyrrel, S. & Cartmell, E. Anaerobic digestion foaming causes – A
1102 review. *Bioresource Technology* **100**, 5546–5554 (2009).
- 1103 72. Duan, J.-L. *et al.* Unraveling anaerobic digestion foaming via association between
1104 bacterial metabolism and variations in microbiota. *ACS EST Eng.* **1**, 978–988
1105 (2021).
- 1106 73. Jiang, C. *et al.* Identification of microorganisms responsible for foam formation in
1107 mesophilic anaerobic digesters treating surplus activated sludge. *Water Research*
1108 **191**, 116779 (2021).
- 1109 74. Bovio-Winkler, P., Cabezas, A. & Etchebehere, C. Database mining to unravel the
1110 ecology of the phylum Chloroflexi in methanogenic full scale bioreactors. *Frontiers*
1111 *in Microbiology* **11**, 603234 (2021).
- 1112 75. McIlroy, S. J. *et al.* Culture-independent analyses reveal novel Anaerolineaceae as
1113 abundant primary fermenters in anaerobic digesters treating waste activated sludge.
1114 *Frontiers in Microbiology* **8**, 1134 (2017).
- 1115 76. Singleton, C. M. *et al.* Connecting structure to function with the recovery of over
1116 1000 high-quality metagenome-assembled genomes from activated sludge using
1117 long-read sequencing. *Nature Communications* **12**, 2009 (2021).
- 1118 77. Jiang, F. *et al.* Recovery of metagenome-assembled microbial genomes from a full-
1119 scale biogas plant of food waste by pacific biosciences high-fidelity sequencing.
1120 *Front. Microbiol.* **13**, (2023).
- 1121 78. Treu, L., Kougias, P. G., Campanaro, S., Bassani, I. & Angelidaki, I. Deeper insight
1122 into the structure of the anaerobic digestion microbial community; The biogas
1123 microbiome database is expanded with 157 new genomes. *Bioresource Technology*
1124 **216**, 260–266 (2016).
- 1125 79. Karst, S. M. *et al.* High-accuracy long-read amplicon sequences using unique
1126 molecular identifiers with Nanopore or PacBio sequencing. *Nat Methods* **18**, 165–
1127 169 (2021).
- 1128 80. Klindworth, A. *et al.* Evaluation of general 16S ribosomal RNA gene PCR primers
1129 for classical and next-generation sequencing-based diversity studies. *Nucleic Acids*
1130 *Research* **41**, e1 (2013).
- 1131 81. Lane, D. J. 16S/23S rRNA sequencing. in *Nucleic Acid Techniques in Bacterial*
1132 *Systematics* (eds. Stackebrandt, E. & Goodfellow, M.) 115–175 (John Wiley and
1133 Sons, Chichester, United Kingdom, 1991). doi:10.1007/s00227-012-2133-0.
- 1134 82. Parada, A. E., Needham, D. M. & Fuhrman, J. A. Every base matters: Assessing
1135 small subunit rRNA primers for marine microbiomes with mock communities, time
1136 series and global field samples. *Environmental Microbiology* **18**, 1403–1414
1137 (2016).
- 1138 83. Caporaso, J. G. *et al.* Ultra-high-throughput microbial community analysis on the
1139 Illumina HiSeq and MiSeq platforms. *The ISME Journal* **6**, 1621–1624 (2012).
- 1140 84. Pinto, A. J. & Raskin, L. PCR biases distort bacterial and archaeal community
1141 structure in pyrosequencing datasets. *PLoS One* **7**, e43093 (2012).

- 1142 85. Apprill, A., McNally, S., Parsons, R. & Weber, L. Minor revision to V4 region SSU
1143 rRNA 806R gene primer greatly increases detection of SAR11 bacterioplankton.
1144 *Aquatic Microbial Ecology* **75**, 129–137 (2015).
- 1145 86. Edgar, R. C. Search and clustering orders of magnitude faster than BLAST.
1146 *Bioinformatics* **26**, 2460–2461 (2010).
- 1147 87. Edgar, R. C. UNOISE2: improved error-correction for Illumina 16S and ITS
1148 amplicon sequencing. *bioRxiv* 81257 (2016) doi:10.1101/081257.
- 1149 88. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing
1150 reads. *EMBnet.journal* **17**, 10–12 (2011).
- 1151 89. Walters, W. A. *et al.* PrimerProspector: De novo design and taxonomic analysis of
1152 barcoded polymerase chain reaction primers. *Bioinformatics* **27**, 1159–1161
1153 (2011).
- 1154 90. R Development Core Team, R. F. F. S. C. & R Core Team. *R: A Language and*
1155 *Environment for Statistical Computing*. Vienna Austria R Foundation for Statistical
1156 *Computing* vol. 1 (R Foundation for Statistical Computing, Vienna, Vienna, 2008).
- 1157 91. RStudio Team. *RStudio: Integrated Development Environment for R*. (RStudio,
1158 PBC., Boston, MA, 2020).
- 1159 92. Oksanen, J. *et al.* *Vegan: Community Ecology Package*. (2019).
- 1160 93. Richard A. Becker, Allan R. Wilks, Ray Brownrigg, Thomas P Minka, & Alex
1161 Deckmyn. maps: Draw Geographical Maps. (2021).
- 1162 94. Dowle, M. & Srinivasan, A. data.table: Extension of `data.frame`. (2019).
- 1163 95. Ogle, D. H., Doll, J. C., Wheeler, A. P. & Dinno, A. FSA: Simple fisheries stock
1164 assessment methods. (2023).
- 1165 96. Mangiafico, S. rcompanion: Functions to support extension education program
1166 evaluation. (2023).
- 1167 97. Pedersen, T. L. Patchwork: the composer of plots. (2020).
- 1168 98. Constantin Ahlmann-Eltze. ggupset: Combination matrix axis for 'ggplot2' to
1169 create 'UpSet' plots. (2020).
- 1170 99. Andersen, K. S. S., Kirkegaard, R. H., Karst, S. M. & Albertsen, M. Ampvis2: an
1171 R package to analyse and visualise 16S rRNA amplicon data. *bioRxiv* 299537
1172 (2018) doi:10.1101/299537.
- 1173 100. Dai, T. *et al.* Identifying the key taxonomic categories that characterize
1174 microbial community diversity using full-scale classification: a case study of
1175 microbial communities in the sediments of Hangzhou Bay. *FEMS Microbiology*
1176 *Ecology* **92**, fiw150 (2016).
- 1177 101. Dueholm, M. K. D. msdueholm/MiDAS5: R-scripts for MiDAS 5: Global
1178 diversity of bacteria and archaea in anaerobic digesters. Zenodo
1179 <https://doi.org/10.5281/zenodo.10982338> (2024).
- 1180 102. Dueholm, M. MiDAS 5: Global diversity of bacteria and archaea in anaerobic
1181 digesters (data for R-scripts). (2024) doi:10.6084/m9.figshare.24219199.v2.
- 1182

1183 **Acknowledgements:**

1184 The project has been funded by the Danish Research Council (grant 6111-00617A,
1185 P.H.N.) and the Villum Foundation (Dark Matter and grant 13351, P.H.N.). We thank
1186 all the involved anaerobic digester plants for providing samples and plant metadata.

1187

1188 **Author Contributions:**

1189 P.H.N. and M.K.D.D. designed the study. M.K.D.D. and P.H.N. wrote the manuscript
1190 and all authors reviewed and approved the final manuscript. M.A., Y.B-F., D.B., C.B.,
1191 M.C.C, Å.D., L.E., C.H., K.K., N.K., C.L., G.L., S.M., V.O., P.O-P., D.P., V.R., M.R.,
1192 J.R., P.E.S., N.T., J.V., J.D.V., C.W. provided samples and metadata. V.R. handled
1193 sampling, DNA extraction and library preparation for DNA sequencing. K.S.A. and
1194 M.K.D.D. performed the bioinformatics analyses. M.K.D.D., A-K.C.P. and K.S.A.
1195 curated metadata and carried out statistical analyses.

1196

1197 **Competing Interests:**

1198 The authors declare no competing interests.