

1 **Disinfection exhibits systematic impacts on the drinking water microbiome.**

2 Zihan Dai<sup>1</sup>, Maria C. Sevellano-Rivera<sup>2</sup>, Szymon T. Calus<sup>1</sup>, Q. Melina Bautista-de los Santos<sup>3</sup>,  
3 A. Murat Eren<sup>4,5</sup>, Paul W.J.J. van der Wielen<sup>6,7</sup>, Umer Z. Ijaz<sup>1</sup>, Ameet J. Pinto<sup>2</sup>.

4

5 <sup>1</sup> Infrastructure and Environmental Research Division, School of Engineering, University of  
6 Glasgow, G12 8LT Glasgow, UK,

7 <sup>2</sup> Department of Civil and Environmental Engineering, Northeastern University, Boston, MA,  
8 USA

9 <sup>3</sup> Department of Civil & Environmental Engineering, University of Michigan, Ann Arbor,  
10 Michigan, USA

11 <sup>4</sup> Department of Medicine, University of Chicago, Chicago, IL, USA

12 <sup>5</sup> Bay Paul Center, Marine Biological Laboratory, Woods Hole, MA, USA

13 <sup>6</sup> KWR Water Research Institute, Nieuwegein, Netherlands.

14 <sup>7</sup> Laboratory of Microbiology, Wageningen University, Wageningen, Netherlands

15 \*Corresponding author: [a.pinto@northeastern.edu](mailto:a.pinto@northeastern.edu)

16

17 **Keywords:** disinfection, drinking water microbiome, selection

## 18 ABSTRACT

19 Limiting microbial growth during drinking water distribution is achieved either by maintaining a  
20 disinfectant residual or through nutrient limitation without the use of a disinfectant. The impact of  
21 these contrasting approaches on the drinking water microbiome is not systematically understood.  
22 We utilized genome-resolved metagenomics to compare the structure, metabolic traits, and  
23 population genomes of drinking water microbiomes across multiple full-scale drinking water  
24 systems utilizing these two-distinct microbial growth control strategies. Microbial communities  
25 cluster together at the structural- and functional potential-level based on the presence or absence  
26 of a disinfectant residual. Disinfectant residual concentrations alone explained 17 and 6.5% of the  
27 variance in structure and functional potential of the drinking water microbiome, respectively,  
28 despite including samples from multiple drinking water systems with variable source waters and  
29 source water communities, treatment strategies, and chemical compositions. The drinking water  
30 microbiome is structurally and functionally less diverse and less variable across disinfected  
31 systems as compared to non-disinfected systems. While bacteria were the most abundant domain,  
32 archaea and eukaryota were more abundant in non-disinfected and disinfected systems,  
33 respectively. Community-level differences in functional potential were driven by enrichment of  
34 genes associated with carbon and nitrogen fixation in non-disinfected systems and  $\gamma$ -aminobutyrate  
35 metabolism in disinfected systems which may be associated with the recycling of amino acids.  
36 Metagenome-assembled genome-level analyses for a subset of phylogenetically related  
37 microorganisms suggests that disinfection may select for microorganisms capable of using fatty  
38 acids, presumably from microbial decay products, via the glyoxylate cycle. Overall, we find that  
39 disinfection exhibits systematic and consistent selective pressures on the drinking water  
40 microbiome and may select for microorganisms able to utilize microbial decay products  
41 originating from disinfection inactivated microorganisms.  
42

## 43 INTRODUCTION

44 Drinking water systems harbor diverse and complex microbial communities in bulk water, biofilms  
45 on pipe wall, suspended solids, and in loose deposits<sup>1-5</sup>. While treatment processes at the drinking  
46 water treatment plants (DWTPs) shape the microbial community that leaves the DWTP<sup>6-9</sup>, multiple  
47 factors can influence the structure and function of the drinking water microbiome in the drinking  
48 water distribution systems (DWDSs). These factors include, but are not limited to, DWDS size,  
49 pipe materials and ages, water age within the DWDS and similar factors within premises plumbing  
50 (PP) in buildings and homes<sup>10-14</sup>. Managing the microbiological quality of drinking water during  
51 transport through the DWDS and into the PP is essential for the provision of safe drinking water.  
52 Unwanted microbial growth and/or changes in the drinking water microbiome composition during  
53 transit through the DWDS and PP are associated with several detrimental outcomes. For instance,  
54 this could lead to proliferation of opportunistic pathogens<sup>15-19</sup> and an eukaryotic microbes<sup>14, 16, 20,</sup>  
55 <sup>21</sup>, taste and odor issues<sup>22</sup>, and impact infrastructure via corrosion damage<sup>23, 24</sup>.

56 Source-to-tap differences in drinking water systems can range from source water type (e.g.,  
57 surface, ground, reuse water), process configurations at the DWTP, heterogeneity and condition  
58 of the DWDS and PP; yet globally there are two fundamental approaches for managing the  
59 drinking water microbiome during transport to the consumer<sup>25</sup>. The first and most widely used  
60 approach involves maintenance of a disinfectant residual (e.g., chlorine) in the DWDS. This is  
61 accomplished by ensuring the water leaving the DWTP has a chlorine residual and/or by using  
62 booster stations in large complex DWDSs to compensate for disinfectant residual decay<sup>26</sup>.  
63 Disinfectant residuals counteract microbial growth through inactivation, thus ensuring stable  
64 microbial concentrations during distribution. While disinfectant residuals are effective in  
65 managing microbial growth in the DWDS, there are some key issues associated with them. These  
66 include aesthetic and corrosion related problems<sup>25, 27, 28</sup>, but more importantly the formation of  
67 harmful disinfection byproducts (DBPs)<sup>29-31</sup>, which are also regulated. Further, there is an  
68 increasing recognition that the disinfectant residuals may be associated with selection of some  
69 opportunistic pathogens<sup>16, 32</sup> and antibiotic resistance genes (ARGs) in drinking water<sup>33-35</sup>.

70 The second approach for managing microbial growth in the DWDS, primarily practiced in parts  
71 of western Europe (e.g., Netherlands, Denmark, and Switzerland), involves distribution of drinking  
72 water without any disinfectant residuals<sup>36</sup>. These systems focus on minimizing nutrient availability

73 in the DWDS to limit microbial growth using high-quality source waters and/or multi-barrier  
74 treatment. While some of these drinking water systems may also use chlorine or other chlorine  
75 compounds (e.g., chlorine dioxide) at the DWTP, they ensure that chlorine is not detectable prior  
76 to distribution. The efficacy of this approach is supported by evidence that incidences of microbial  
77 contamination and associated waterborne illnesses are comparable to systems that maintain a  
78 disinfectant residual<sup>25, 37</sup>. This suggest that with appropriate source water quality management,  
79 treatment, and well maintained infrastructure, drinking water can be safely distributed without  
80 disinfectant residuals<sup>25</sup>.

81 Despite reports of comparable biological water quality between systems with and without  
82 disinfectant residuals, there are a limited number of studies that have systematically compared the  
83 microbial community between these two types of systems. Bautista et al (2016)<sup>38</sup> conducted a  
84 meta-analyses study involving collation, curation, and comparison of 16S rRNA gene amplicon  
85 sequencing data from previously published datasets. While this study was confounded by  
86 methodological differences between datasets being used, the key conclusions were that  
87 presence/absence of disinfectant residuals impact microbial community structure and membership  
88 and that systems without disinfectant residuals are more diverse than their disinfected counterparts.  
89 Recently, Waak et al (2019)<sup>39</sup> compared biofilms between two drinking water systems, one  
90 chloraminated systems and one without a disinfectant residual. Consistent with previous findings  
91 they observed higher cell numbers and higher diversity in the system without disinfectant residual,  
92 with higher proportional abundance (proportion of total community) of deleterious microbes (i.e.,  
93 mycobacteria, nitrifiers, corrosion causing bacteria) in the chloraminated system. Both, Bautista  
94 et al (2016)<sup>38</sup> and Waak et al (2019)<sup>39</sup> utilized gene-targeted assays (i.e., 16S rRNA gene) to probe  
95 drinking water microbiome composition and its differences. While gene-targeted assays can  
96 provide valuable information on microbial community structure and membership information,  
97 they do not provide insight into metabolic differences that may drive the observed differences in  
98 community structure. Further, gene-targeted assays can be limited by primer-bias and can result in  
99 non-detection of microbial community members. Both challenges can be overcome by utilizing  
100 metagenomics which can provide insights into structure and functional potential of a microbial  
101 communities without being biased against or towards specific community members. This comes  
102 with the limitation that differences between samples/systems emerging from low-abundance  
103 microbes may not be detected as this may require ultra-deep sequencing.

104 We used metagenome analyses and genome-resolved metagenomics to investigate the potential  
105 influence of disinfectant residuals on the drinking water microbiomes by comparing drinking water  
106 systems from the United Kingdom (with disinfectant residual) and the Netherlands (without  
107 disinfectant residual). The goals of this study were (1) to determine the extent to which disinfectant  
108 residual shapes the structure and functional potential of the drinking water microbiome, (2) to  
109 determine whether the selective pressures of disinfection are conserved across drinking water  
110 systems, and (3) identify metabolic pathways underpinning differences in structure and functional  
111 potential of the drinking water microbiome. Addressing these questions across different drinking  
112 water systems with inherent system-to-system variability (e.g., source water, water chemistry,  
113 treatment process, etc.) but one consistent difference - i.e., presence or absence of disinfectant  
114 residual - will help highlight disinfection that are conserved and thus generalizable across systems.

## 115 **MATERIALS AND METHODS**

### 116 **Sample collection and processing.**

117 Drinking water samples were collected from 12 drinking water systems in Netherlands (n=5)  
118 between October to December 2013 (Non-disinfected, i.e. ND) and the United Kingdom (n=7)  
119 between April to August in 2015 (Disinfected, i.e. D). Samples were collected at two to four  
120 locations in each DWDS which resulted in 23 D and 18 ND samples. A total 15 liters of water was  
121 filtered through three sterile Sterivex filters with 0.22 $\mu$ m pore size polyethersulfone membrane  
122 (EMD Millipore<sup>TM</sup> SVGP01050) using a peristaltic pump (Watson-Marlow 323S/D) to harvest  
123 microbial cells. Immediately after filtration, the membranes were removed aseptically from the  
124 Sterivex cartridge, cut into pieces and then transferred to Lysing Matrix E tubes. The membranes  
125 were stored at 4°C for 24 hours or less before being transported to the laboratory and stored at -  
126 80°C. Further details of sample treatments and preservation are described in Sevillano-Rivera et  
127 al.<sup>35</sup>, along with detailed description of chemical analyses. Briefly, Orion 5 Star Meter (Thermo  
128 Fisher Scientific, Waltham, MA) was used to measure temperature, pH, conductivity and dissolved  
129 oxygen, total chlorine, and phosphate was also determine on-site using DR 2800 VIS  
130 Spectrophotometer (Hach Lange, the UK) and EPA approved HACH kits. Nitrogen species were  
131 measured according to standard method, 4500-NH3-F for ammonia, 4500-NO2-B for nitrite, and

132 4500-NO3-B for nitrate respectively in laboratory<sup>40</sup>, while total organic carbon (TOC) was  
133 determined using Shimadzu TOC-LCPH Analyzer (Shimadzu, Kyoto, Japan).

#### 134 **DNA extractions.**

135 The total genomic DNA was extracted directly from filter membranes using Maxwell16 DNA  
136 extraction system (Promega) and LEV DNA kit (AS1290, Promega, Madison, WI, US). The filters  
137 with collected biomass in lysing matrix E tubes were incubated with 300 $\mu$ L of lysing buffer and  
138 30 $\mu$ L of Proteinase K and incubated at 56°C. A total of 500 $\mu$ L of chloroform:isoamyl alcohol  
139 (24:1, pH 8.0) was added to the tube, vortexed and this was followed by bead beating for 40 s at 6  
140 m/s using a FastPrep 24 instrument (MP Biomedicals, Santa Ana, CA, USA), and centrifugation  
141 at 14,000g for 10 min. The bead beating and centrifugation steps were repeated twice more with  
142 transfer of supernatant to clean tube followed by replacement of the aqueous phase with fresh  
143 lysing buffer. The aqueous phase was then subject to DNA purification using the Maxwell LEV  
144 DNA kit. The extracted DNA was quantified using Qubit HS dsDNA assay with Qubit 2.0  
145 Fluorometer (Life Technologies, UK). Negative controls consisting of reagent blanks (no input  
146 material) and filter blanks (filter membranes from unused Sterivex filters) were processed  
147 identically as the samples for DNA extraction. Genomic DNA extracted from mock community,  
148 consisting of 10 organisms, detailed previously<sup>35</sup>, was spiked into negative controls extracted  
149 (n=8) from the reagent and filter blanks. These negative controls were also included in following  
150 library preparation and high-throughput sequencing (see below).

#### 151 **Library preparation and Illumina sequencing.**

152 Sequencing libraries were prepared using the Nextera XT DNA Sample Preparation Kit (Illumina  
153 Inc.). All DNA extracts (including negative controls) were cleaned up with HighPrep PCR  
154 magnetic beads (MagBio Inc.) to remove short fragments after library preparation and quantified  
155 with qPCR according to Illumina guidelines. All libraries were pooled together in equimolar  
156 proportion and pooled library was quantified with Qubit HS dsDNA assay and further concentrated  
157 using HighPrep PCR magnetic beads (MagBio Inc). Metagenomic sequencing on prepared  
158 libraries were performed on four lanes of Illumina HiSEQ 2500 flow cell (2x250-bp read length,  
159 Rapid Run Mode) at University of Liverpool Centre for Genomic Research (Liverpool, UK).

## 160 **Metagenomic read based analyses.**

161 The FASTQ files were trimmed using Cutadapt v1.2.1 (Martin 2014) with a '-O 3' flag, and Sickle  
162 v1.200 (Joshi and Fass 2011) using a threshold of window quality score ( $\geq 20$ ) and read length  
163 after trimming ( $\geq 10$  bp). A further trimming was applied using Trimmomatic v0.35<sup>41</sup> to remove  
164 any remaining Illumina Nextera adaptors and trim reads according to quality score with a 4-base  
165 wide sliding window and a minimum average quality score of 20 and singlet reads were excluded  
166 in downstream analyses as well. To estimate metagenome diversity and coverage for each sample,  
167 Nonpareil 3.0<sup>42</sup> was used in kmer mode on the quality filtered reads. Diversity and coverage  
168 information for each metagenome was estimated using command 'Nonpareil.set()' in R package  
169 'Nonpareil'. MicrobeCensus<sup>43</sup> was used on quality trimmed reads to estimate average genome size  
170 across samples with flag '-n 100000000' for all samples. To eliminate the potential effects of  
171 bacteria with small genomes (i.e., *Patescibacteria*) on average genome size estimations, pre-  
172 processed reads were mapped against 12 *Patescibacteria* metagenome assembled genomes  
173 (MAGs) from this study (see below) and 1,037 *Patescibacteria* genomes from GTDB-tk<sup>44</sup>. The  
174 reads mapped in proper pair to *Patescibacteria* were removed using samtools ('-F2' flag).  
175 MicrobeCensus was used again to estimate average genome size using the same parameters.

## 176 **Metagenome assembly and mapping.**

177 Filtered pair-ended reads were then pooled from each drinking water system for co-assembly,  
178 which resulted in 12 paired-end FASTQ files for co-assembly, including seven from disinfected  
179 (Dis) and five from non-disinfected (NonDis) systems. *De novo* co-assembly was performed using  
180 MetaSPAdes v3.10.1<sup>45</sup> with recommended k-values for 2x250bp reads (21,33,55,77,99,127). All  
181 scaffolds shorter than 500bp were discarded and UniVec\_Core build 10.0 (National Center for  
182 Biotechnology Information 2016) was used for contamination vector screening and any scaffold  
183 with a significant hit to the UniVec database was removed. Reads from each samples were then  
184 mapped back to the filtered scaffolds using BWA-MEM v0.7.12 with default settings<sup>46</sup>.

185 To eliminate the scaffolds that may have originated from sample or post-processing contamination,  
186 reads from negative controls were first mapped back to mock community genomes using BWA-  
187 MEM v0.7.12<sup>46</sup>, and all reads not mapped in proper pair were extracted using samtools v1.3.1 (Li  
188 et al. 2009) with '-f2' flag and were considered "contaminant reads". Sample reads (S),

189 contaminant reads (C) and negative control reads (NC) were mapped back to filtered scaffolds in  
190 each co-assembly. Properly-paired mapped reads were extracted using samtools v1.3.1 with '-f2'  
191 flag from the BAM files. Relative abundance and normalized coverage deviation of each scaffold  
192 was calculated using reads from samples and those identified as contaminant reads in negative  
193 controls:

$$194 \quad \text{Relative abundance}_S = \frac{\text{Scaffold coverage}_S}{\sum_{i=1}^n \text{Scaffold coverage}_S}$$

$$195 \quad \text{Relative abundance}_C = \frac{\text{Scaffold coverage}_C}{\sum_{i=1}^n \text{Scaffold coverage}_{NC}}$$

$$196 \quad \text{Normalized coverage deviation} = \frac{\text{Standard deviation of scaffold coverage}}{\text{Average scaffold coverage}}$$

197 To distinguish true scaffolds from contamination, relative abundance (RA) and normalized  
198 coverage deviation (NCD) estimated using sample reads (S) and contaminant reads (C) was  
199 compared for all scaffolds:

$$200 \quad \text{Scaffold} = \begin{cases} \text{True scaffold,} & \text{if } \begin{matrix} RA_C = 0 \\ RA_S > RA_C \text{ and } NCD_S < NCD_C \end{matrix} \\ \text{Contaminant scaffold,} & \text{if } \begin{matrix} RA_S = 0 \\ RA_C > RA_S \text{ and } NCD_C < NCD_S \end{matrix} \end{cases}$$

201 True scaffolds, the scaffolds with higher RA and lower NCD in samples compared to negative  
202 controls, were kept for downstream analyses while contaminant scaffolds were excluded from all  
203 further analyses.

#### 204 **Nucleotide and protein composition analyses.**

205 MASH v1.1.1<sup>47</sup> was used to estimate the dissimilarity between samples using quality filtered reads  
206 (with '-r' and '-m 2' flags) and dissimilarity between drinking water systems using true scaffolds  
207 with the sketch size of 100000. Prodigal v2.6.3<sup>48</sup> was used to identify open reading frames (ORFs)  
208 in the true scaffolds and translate ORFs to protein-coding amino acid sequences. Following  
209 prediction and translation, HMMER v3.1b2<sup>49</sup> was used to annotate ORFs against the Pfam  
210 database v31.0<sup>50</sup> with a maximum e-value of  $1e - 5$  and curated bit score thresholds (the gathering  
211 thresholds). Subsequently, MASH distances were calculated between drinking water



212 metagenomes using predicted ORFs, as well as Pfam annotated proteins with the sketch size of  
213 100000 and '-a' flag.

#### 214 **Taxonomic classification and phylogenetic analyses.**

215 The program 'cmsearch' was implemented in Infernal v1.1.2<sup>51</sup> to search scaffolds against SSU  
216 rRNA covariance models (CMs) for bacteria, archaea and eukaryota; these are default models used  
217 by SSU-ALIGN v0.1<sup>52</sup> using HMM-only approach and only significant hits were considered. The  
218 results were filtered according to length ( $\geq 100$  bp alignment) and e-value ( $< 1e - 5$ ). SSU rRNA  
219 sequences detected in contaminant scaffolds were removed and if more than one SSU gene  
220 sequence was found on a single scaffold, only the longest SSU gene sequence was retained.  
221 Relative abundance of each SSU gene sequence was calculated for each sampling location as  
222 follows:

$$223 \quad \text{RPKM}_{SSU}^i = \frac{\text{Scaffold coverage}^i}{\sum_{i=1}^n \text{SSU containing Scaffold coverage per Mb}^i \times \text{Scaffold length per kb}^i}$$

$$224 \quad \text{Relative abundance}_{SSU}^i = \frac{\text{RPKM}_{SSU}^i}{\sum_{i=1}^n \text{RPKM of scaffold containing SSU gene}^i}$$

225 SSU rRNA gene sequences were classified using Mothur v1.33.3 (Schloss et al. 2009) with SILVA  
226 database<sup>53</sup> (Release 132) with a minimum confidence threshold of 80%.

#### 227 **Annotation and Comparison of functional orthologies and modules between samples**

228 The protein-coding sequences were searched against KOfam, a HMM profile database for KEGG  
229 orthology<sup>54</sup> with predefined score thresholds using KofamScan<sup>55</sup>. Only KEGG orthologies (KO)  
230 identified on scaffolds with ( $> 1x$ ) coverage for each sample and those detected more than once  
231 across samples within a single drinking water system were retained for further analyses. Average  
232 read count for each KO was calculated using scaffold coverage, average length of reads mapped,  
233 and total number of reads mapped to each scaffold in a sample using above equations. To assess  
234 functions at KEGG module level, BRITE hierarchy file was retrieved from KEGG website, and  
235 KO's were categorized into KEGG modules. The abundance of KEGG module in each sample was  
236 calculated using the median abundance of the detected KEGG orthologies within each module.  
237 The completeness of each KEGG module was calculated using 'KO2MODULEclusters2.py'.

## 238 **Metagenome binning and refining.**

239 Anvi'o (versions: v2.2.2, v2.4.0, v4 and v5.1)<sup>56</sup> was used for metagenome binning and refining.  
240 Briefly, CONCOCT<sup>57</sup> integrated in Anvi'o was used to cluster scaffolds (longer than 2500 bp) into  
241 metagenome bins using tetra-nucleotide composition and coverage information across all samples  
242 within each metagenomic co-assembly. The 'merge' method of CheckM v1.0.7<sup>58</sup> was used to  
243 identify the bins that that may emerge from the same microbial population, but may have been  
244 separated during automated binning process. Following merging of compatible bins, RefineM  
245 v0.0.21<sup>59</sup> was used to automatically refine bins according to genomic properties (i.e., the mean GC  
246 content, tetra-nucleotide signature and coverage) and taxonomic classification. The completeness  
247 and redundancy of each refined bin was estimated using CheckM based on collections of lineage  
248 specific single-copy genes resulting in a total of 154 bins with greater than >50% completeness.  
249 Among these bins, 130 bins had a redundancy of <10% redundancy, while 24 bins are with >10%  
250 redundancy. Further manual curation of these bins was performed using Anvi'o, resulting in 156  
251 curated metagenome assembled genomes (MAGs). The 156 MAGs were de-replicated using dRep  
252 v2.2.2<sup>60</sup> and MAGS with >10% redundancy were discarded which resulted in 115 dereplicated  
253 MAGs with completeness >50% and reduncancy <10%. All raw sequencing data and dereplicated  
254 MAGs are available on NCBI at BioProject number PRJNA533545.

## 255 **MAG-level analyses**

256 Taxonomy assignment of MAGs was performed using GTDB-Tk v0.1.3<sup>44</sup> with the flag  
257 'classify\_wf'. Genome sizes of MAGs were estimated by multiplying the number of nucleotides  
258 in the MAG with the inverse of the CheckM estimated completeness. The MAGs were annotated  
259 using the HMM profile database for KEGG orthology with predefined score thresholds using  
260 KofamScan<sup>55</sup>. The KO's for each MAG were then categorized into modules based on BRITE  
261 hierarchy file retrieved from KEGG<sup>54</sup>, and the completeness of KEGG modules in each genome  
262 was calculated using script 'KO2MODULEclusters2.py'. Anvi'o was used to extract a collection  
263 of 48 single-copy ribosomal proteins<sup>61</sup> from each MAG using 'anvi-get-sequences-for-hmm-hits'  
264 with a maximum number of missing ribosomal proteins of 40. Subsequently, a phylogenetic tree  
265 was reconstructed using concatenated alignment of ribosomal proteins sequences using FastTree  
266 v2.1.7<sup>62</sup>. Interactive Tree Of Life (iTOL) v4 (Letunic and Bork 2007) was used to visualize the  
267 phylogenetic tree.

268 Program 'Union' in EMBOSS v6.6.0.0<sup>63</sup> was used to concatenate all scaffolds in each MAG into  
269 a single sequence. Reads from all samples were cross-mapped to all MAGs using BWA-MEM  
270 v0.7.12 with default settings and proportion of each nucleotides in MAG covered by at least 1x  
271 coverage was determined using BEDtools<sup>64</sup>. A MAG was considered detected in a sample if  $\geq 25\%$   
272 of its bases were covered by at least one read from the corresponding sample. This approach was  
273 used to determine whether MAGs were detected in all the samples. Further, the MAGs were binned  
274 into four categories based on their detection/non-detection within samples. Specifically, MAGs  
275 were divided into "D-only" if there were detected in  $\geq 20\%$  of the samples from the disinfected  
276 systems and not detected in any samples from the non-disinfected systems, "ND-only" if there  
277 were detected in  $\geq 20\%$  of the samples from the non-disinfected systems and not detected in any  
278 samples from the disinfected systems, "both" if there were detected in  $\geq 20\%$  of disinfected and  
279 non-disinfected systems, while the remaining MAGs were classified in the "other" category.  
280 Subsequently, reads from all samples were cross-mapped back to all the MAGs using BMap  
281 v38.24<sup>65</sup> with a minimum identity of 90%, and 'ambiguous=best' and 'pairedonly=t' flags. After  
282 filtering for detection (see above), reads per kilobase of per million reads (RPKM) for each MAG  
283 and each sample were calculated using equation:

$$284 \quad \text{RPKM} = \frac{\text{Number of reads mapped to MAG}}{\text{Total number of reads in sample per Million} \times \text{MAG length in kbp}}$$

## 285 **Statistics**

286 Differences between disinfected and non-disinfected systems for (1) Mash distances distributions,  
287 (2) relative abundances were determined using Permutational ANOVA and Pearson's correlations  
288 between pairwise mash distances were estimated in R. BioEnv in "sinkr"  
289 (<https://github.com/menugget/sinkr>) and "vegan"<sup>66</sup> packages were used to identify  
290 environmental parameters (i.e., water chemistry) and their combinations that explain the  
291 differences in the structure (i.e., Mash distances between samples estimated using reads) and  
292 functional potential (i.e., Bray Curtis distance estimated between samples using KO abundance  
293 (i.e., RPKM). BioEnv permutes through  $2^{n-1}$  possible combination of selected environmental  
294 parameters, 511 combinations in this case, and selects the combinations of scaled environmental  
295 variables which capture maximum correlation between dissimilarities of community datasets water  
296 chemistry and microbial community structure or functional potential. While, BioEnv analyses

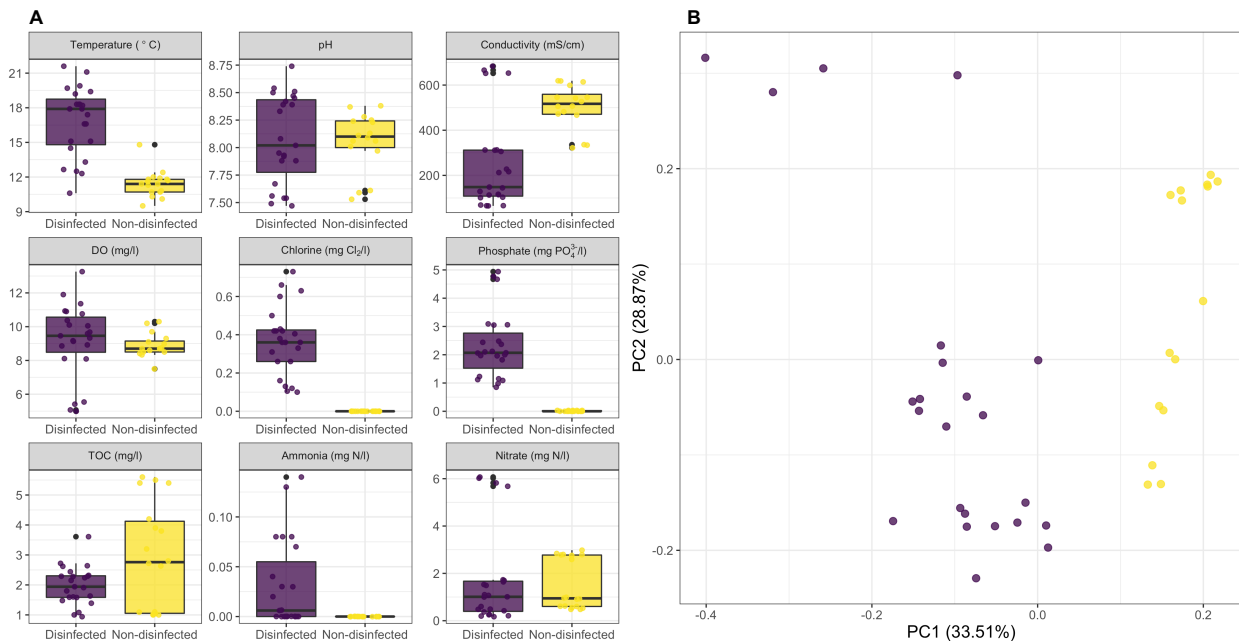
297 identified combination of variables that are highly correlated with differences in microbial  
298 community structure of functional potential, it does not identify the proportion of variance in  
299 microbial community structure of functional potential explained by individual variables or their  
300 combination. To this end, we used distance-based redundancy analysis (dbRDA) to perform  
301 constrained ordinations on community structure and functional potential to bypass the limitation  
302 of usual RDA and CCA, which can only use Euclidean distance measure. Function `dbrda()` from  
303 'vegan' package was used with pairwise Mash distances calculated between samples estimated  
304 using reads based Mash distance and Bray-Curtis distances based on KO RPKM to investigate  
305 relationships between the environmental variables and community data on both nucleotide  
306 composition and KO level. The function `varpart()` in the vegan package was used to determine the  
307 fraction of variation captured parameters identified as significantly associated with read-based  
308 Mash and KO relative abundance-based Bray-Curtis distance matrices. DESeq2 package v1.18.1<sup>67</sup>  
309 was used to identify differentially abundant KEGG modules between disinfected and non-  
310 disinfected systems by only considering KEGG modules with a maximum of one block missing  
311 and equal to or greater than 50% complete. The median scaffold-length normalized read count of  
312 KO's within each module were used in DESeq2 analyses with a maximum adjusted *P*-value of  
313 0.005.

## 314 RESULTS AND DISCUSSION

### 315 **Water quality parameters across disinfected and non-disinfected DWDS.**

316 Sampling was conducted in seven DWDSs with disinfectant residual between April-August of  
317 2013 and at five DWDSs without disinfectant residual between October-December 2015. The  
318 water chemistry varied between the DWDSs considering they were supplied by different DWTPs,  
319 our sampling campaign also captures seasonal differences between locations (Figure 1) (Table S1).  
320 Specifically, water temperatures were higher (~5°C) for the disinfected samples compared to the  
321 non-disinfected samples. While the pH, DO, nitrogen species (i.e., ammonium and nitrate) and  
322 TOC measurements were not significantly different between disinfected and non-disinfected  
323 samples, the measured phosphate and total chlorine concentrations were significantly different  
324 ( $p < 0.05$ ). Specifically, the average total chlorine concentrations in disinfected systems 0.37 mg  
325  $\text{Cl}_2/\text{l}$  (range: 0.1-0.73 mg  $\text{Cl}_2/\text{l}$ ) while no disinfectants residuals were measurable in the non-  
326 disinfected systems. The average phosphate concentrations were 2.3 mg  $\text{PO}_4^{3-}/\text{l}$  while no

327 phosphate was measurable in non-disinfected samples. Phosphate was higher in the disinfected  
328 systems as it is likely to be used for corrosion control<sup>68</sup>. While we were unable to obtain  
329 information on source water type (i.e., ground vs surface water) used for production of drinking  
330 water supplied to the sampled DWDS, conductivity measurements suggested DWDS in both  
331 systems were supplied by a DWTPs drawing from surface and ground water sources (Figure 1).



**Figure 1:** Summary of water chemistry parameters measured for samples collected from disinfected (purple) and non-disinfected systems (yellow). (B) Principle component analyses using Euclidean distances for measured water chemistry parameters indicates distinct clustering of samples from disinfected and non-disinfected systems.

332

### 333 **Summary of metagenomic data set.**

334 Metagenomic analyses was used to assess the association between presence/absence of disinfectant  
335 residual with the structure and functional potential of the drinking water microbiome. A total of  
336 41 drinking water samples were collected from DWDSs with (i.e., chlorine) from the United  
337 Kingdom (n=23), while those collected from the Netherlands (n=18) did not have a disinfectant  
338 residual. Quality trimming of raw metagenomic data resulted in the retention of 638 million paired-  
339 end reads. Co-assembly for each drinking water system was carried out by combining reads from  
340 individual sampling location within each drinking water system (Table 1). *De novo* co-assembly  
341 generated 0.04-1.81 million true scaffolds for each sampling location after discarding scaffolds  
342 shorter than 500bp and contaminant scaffolds (Table 1) with an N50 value ranged from 775 bp to

343 3300 bp. The proportion of quality trimmed reads mapping back to true scaffolds ranged from 67%  
344 to 99% (Table 1) across all samples.

345 **Table 1:** Sequencing and de novo co-assembly statistics for metagenomes from 12 drinking  
346 water systems.

Drinking water system	Paired end Reads (millions)	Scaffolds (>500 bp)	True scaffolds	True scaffold assembly size (Mbp)	% Mapped reads	GC content (%)	N50 (bp)	ORFs per Mbp	Coding density
D1	195.73	555493	546375	615.10	99.02	54.66	1131	1403.68	0.48
D2	46.87	38567	36733	53.03	96.24	55.34	2112	1419.84	0.64
D3	17.40	192457	190882	249.69	91.15	57.82	1531	1498.24	0.63
D4	36.01	123852	122486	204.78	93.03	57.57	3300	1316.54	0.60
D5	36.74	227196	225149	269.12	88.73	59.09	1313	1527.12	0.60
D6	17.39	42209	41459	57.23	95.89	59.16	1641	1504.23	0.65
D8	19.4	77973	76996	108.07	95.38	61.07	1751	1475.71	0.68
ND1	45.52	521371	517773	472.02	83.82	53.75	855	1803.21	0.61
ND2	25.98	363819	361304	316.18	75.03	53.44	802	1807.05	0.56
ND3	48.63	667992	663968	562.73	81.63	52.93	775	1838.06	0.60
ND4	17.78	164328	163361	143.22	66.73	56.48	808	1822.84	0.63
ND5	130.92	1812573	1804048	1834.75	93.74	56.38	1005	1672.04	0.60

347 D=disinfected, ND=non-disinfected, N50=minimum contig length that account for 50% of the  
348 bases, ORF=open reading frame.

### 349 **Non-disinfected systems are more diverse than disinfected systems.**

350 Non-disinfected systems were significantly ( $p < 0.0001$ ) more diverse compared to systems that  
351 maintained a disinfectant residual (Figure 2A) based on Nonpareil estimated diversity index<sup>42</sup>. This  
352 observation is consistent with previous comparisons of bulk water<sup>69</sup> and biofilm<sup>39</sup> samples from  
353 disinfected and non-disinfected systems. Lower diversity in disinfected systems is likely due to  
354 stronger selective pressure of the disinfectant residual as compared to that nutrient limitation in  
355 non-disinfected systems. As a result of the higher diversity in non-disinfected systems, the  
356 metagenomic sequencing for these samples provided significantly lower coverage of the sampled  
357 microbial community (Figure 2B) as compared to systems with a disinfectant residual ( $p <$   
358 0.0001).

359 **Microbial community membership and structure is different between disinfected and non-**  
360 **disinfected systems.**

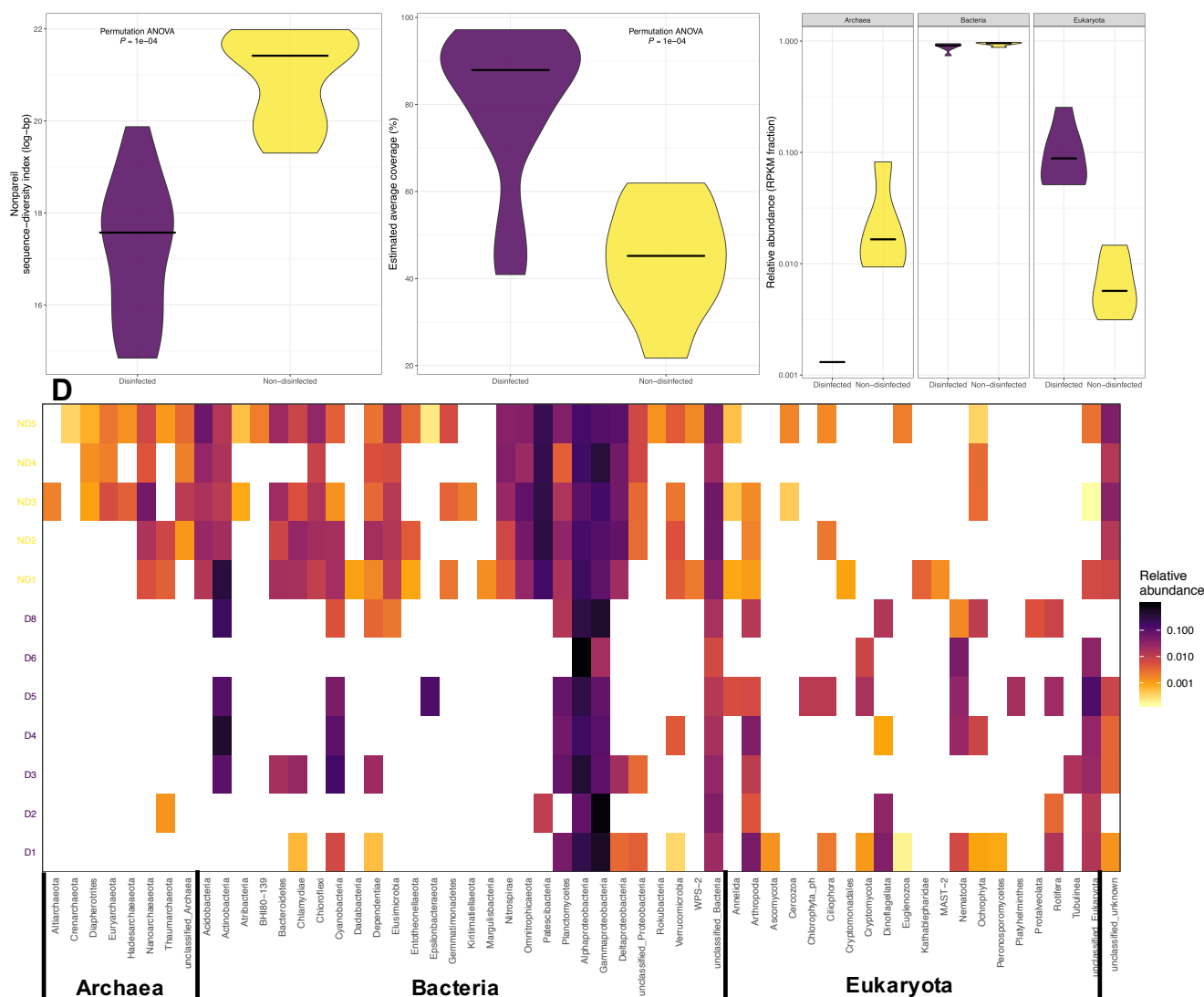
361 We used 2,872 small-subunit (SSU) rRNA genes (2742 genes > 100 bp) identified on the  
362 assembled scaffolds to determine community membership and structure across sampling locations  
363 (Supplementary file 1, Supplementary table 2). While bacteria were dominant members of the  
364 drinking water microbiome in both types of systems (2C, 2D), the relative abundance of archaea  
365 and eukaryota were dependent on the presence/absence of disinfectant residual (Figure 2C, 2E).  
366 Specifically, the relative abundance of eukaryota was higher in disinfected systems as compared  
367 to non-disinfected system (2C), while archaea were ubiquitous across non-disinfected samples  
368 (Figure 2C, E) they were only detected in a single disinfected sample (D2). Non-disinfected  
369 systems were taxonomically more diverse, with respect to bacteria and archaea, as compared to  
370 disinfected systems. Specifically, a total of 14 bacterial and 6 archaeal phyla were detected in one  
371 or more non-disinfected systems that were not detected in any of the disinfected systems. Several  
372 of these unique phyla, while not dominant in non-disinfected systems, were detected at relative  
373 abundances between 1-5% (e.g., *Nitrospirae*, *Nanoarchaeota*).

374 The bacterial community was dominated by *Proteobacteria*, in particular *Alphaproteobacteria* and  
375 *Gammaproteobacteria*, in both disinfected and non-disinfected systems with *Deltaproteobacteria*  
376 being much more prevalent and abundant in non-disinfected systems (Figure 2D). *Actinobacteria*  
377 were more abundant than *Proteobacteria* in two drinking water systems and constituted 44% and  
378 33% of the community in systems D4 and ND1, respectively. Overall, the relative abundance of  
379 *Proteobacteria* was higher in disinfected systems, ranging from 28% to 90%, as compared to non-  
380 disinfected systems, ranging from 30% to 57%. *Patescibacteria* was the second most abundant  
381 phylum across non-disinfected systems, constituting 15% to 29% of the SSU rRNA genes, while  
382 they were only detected in one disinfected sample (D2) with a relative abundance of 1%. Within  
383 *Patescibacteria*, *Parcubacteria* were the most commonly detected phyla followed by  
384 *Microgenomatia* and *Gracilibacteria*.

385 The observed differences between disinfected and non-disinfected DWDS for bacteria and archaea  
386 are largely consistent with a previous meta-analysis of amplicon sequencing data from the 16S  
387 rRNA gene<sup>69</sup>. In contrast to bacteria and archaea, results from eukaryotes, which have not been  
388 systematically investigated in the drinking water microbiome, were surprising in terms of their

389 higher relative abundance eukaryotic in disinfected as compared to non-disinfected systems. For  
390 instance, SSU rRNA genes associated *Nematoda* were detected in nearly every disinfected system,  
391 but were not detected in non-disinfected systems. Specifically, SSU rRNA genes from two free-  
392 living nematode genera, i.e. *Araeolaimida* and *Monhysterida*, were detected in five of the eight  
393 disinfected systems. Similarly, SSU rRNA genes from the phylum *Rotifera* were only detected in  
394 disinfected systems and were largely associated with the monogonont rotifers within the genera  
395 *Ploimida*. While the relative abundance of scaffolds determined to be of eukaryotic origin was  
396 higher in disinfected compared to non-disinfected systems, this does not mean that eukaryotes  
397 were proportionally larger part of the drinking water microbiome in disinfected compared to the  
398 non-disinfected systems. Genome sizes of picoeukaryotic microbes can be orders of magnitude  
399 larger than that of bacteria and archaea and vary significantly between picoeukaryotes themselves.  
400 Further, the higher overall diversity and lower sequencing coverage (Figure 1) could also have  
401 resulted in under sampling of the eukaryotic community in non-disinfected systems.





**Figure 2:** Comparison of (A) diversity and (B) coverage between disinfected and non-disinfected drinking water systems estimated using Nonpareil. (C) Comparison of relative abundance of bacterial, archaeal, and eukaryotic communities in drinking water systems with and without disinfectant residuals. (D) Log<sub>10</sub> transformed relative abundance of different phyla (classes for phylum Proteobacteria) across sampling location for the bacteria, archaea, and eukaryota.

402

403 **Drinking water systems cluster at the nucleotide level based on presence/absence of**  
 404 **disinfectant residuals.** Samples (for read based analyses) and drinking water systems (for scaffold  
 405 based analyses) clustered with each other based on the presence/absence of disinfectant residuals  
 406 (Figure 3A and 3B) based on Mash distance estimates. We further evaluated the significance and  
 407 explanatory power of measured water chemistry parameters in explaining the observed clustering  
 408 between disinfected and non-disinfected systems. To do this, we initially performed BioEnv

409 analyses to identify water chemistry parameters and their combinations that were highly correlated  
410 with observed Mash distances between samples (Supplementary Table 3). This identified chlorine  
411 as being strongly correlated with the Mash distances between samples ( $R=0.54$ ,  $p<0.001$ ) while  
412 the maximum correlation between water chemistry and Mash distances was observed for a  
413 combination of chlorine, phosphate, and TOC ( $R=0.62$ ,  $p<0.001$ ). We subsequently utilized  
414 dbRDA to independently determine the environmental/water chemistry variables most  
415 significantly associated with Mash distances between samples. While chlorine was identified as a  
416 significant variable ( $p<0.01$ ), dbRDA identified conductivity ( $p<0.001$ ) and DO ( $p<0.01$ ) as  
417 significant variables. Finally, variance partitioning analyses was used to determine the proportion  
418 of variance in the Mash distance matrices explained by individual and combination of variables  
419 identified as significant by dbRDA (Table S5). This resulted in chlorine, conductivity, and DO  
420 individually explaining  $\sim 17\%$ ,  $12\%$ , and  $1\%$  of the variance in the Mash distance matrix, with  
421  $\sim 60\%$  of the variance unexplained by these three variables.

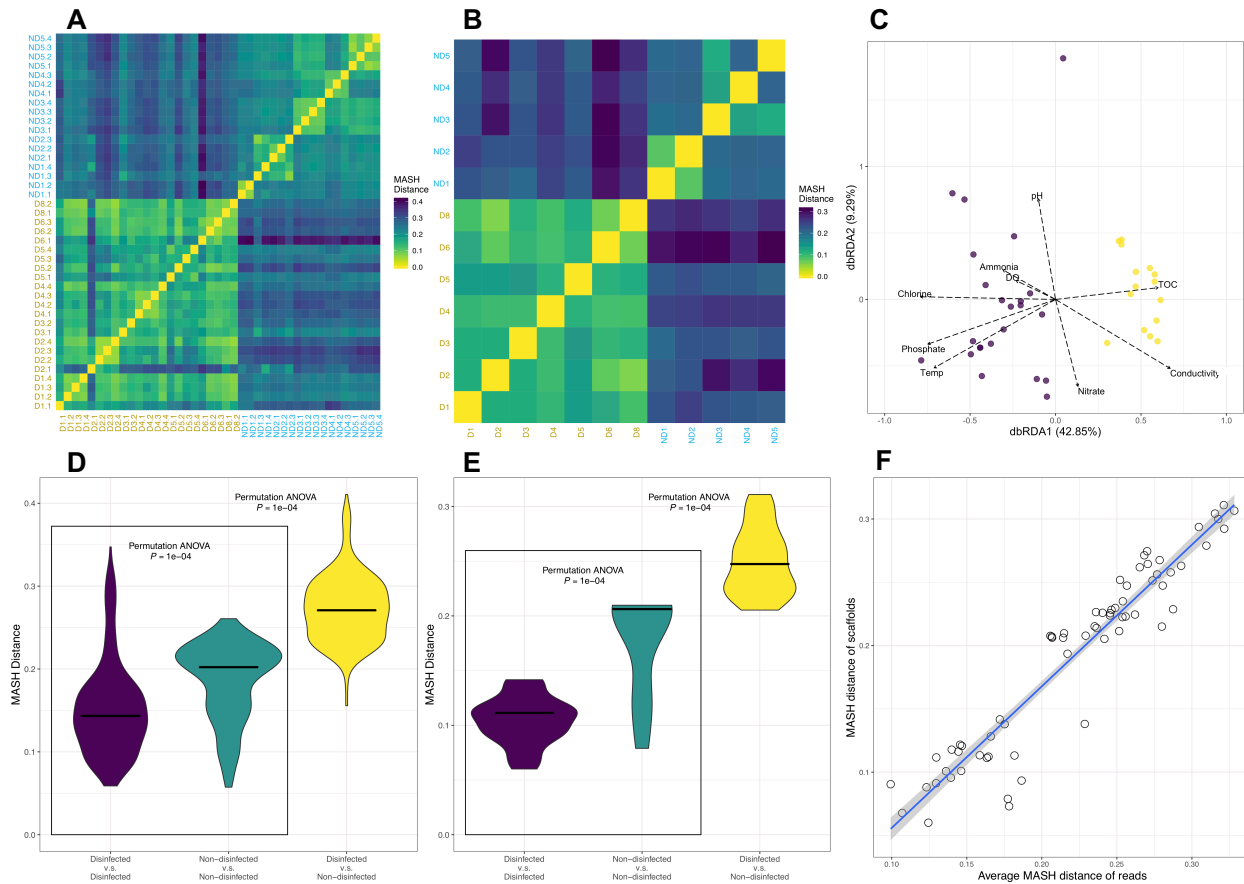
422 We further compared the distribution of Mash distances between drinking water metagenomes  
423 within disinfected, within non-disinfected, and between disinfected and non-disinfected systems.  
424 Mash distances between drinking water metagenomes from disinfected systems were significantly  
425 different ( $p < 0.0001$ ) and exhibited a lower mean for disinfected as compared to non-disinfected  
426 systems. Further, the pairwise Mash-distances between disinfected and non-disinfected systems  
427 were significantly different and higher from those estimated within each category (i.e., disinfected  
428 or non-disinfected). This was consistent for both read- and scaffold-based analyses (Figure 3D,  
429 3E). Finally, the average pairwise Mash distances estimated using reads (i.e., between samples)  
430 and scaffolds (i.e., between DWDSs) were highly correlated (Pearson's  $R = 0.95$ ,  $P < 0.05$ ) (Figure  
431 3C), indicating the *de novo* assembly process did not result in loss of information on factors driving  
432 the differences between disinfected and non-disinfected systems.

433 These analyses provide a few key insights. First, Mash distance-based (both read and scaffold  
434 based) clustering of samples occurs depending on presence and absence of disinfectant residual  
435 suggests that the microbial communities are more similar within each group (i.e., disinfected and  
436 non-disinfected) and dissimilar between the two groups (i.e., disinfected vs non-disinfected).  
437 Second, while disinfected and non-disinfected samples cluster distinctly from each other,  
438 disinfected systems exhibit lower nucleotide-level heterogeneity as compared to their non-

439 disinfected systems indicating that the factors governing microbial community in disinfected  
440 systems likely impose stronger selective pressures on the microbial community as compared to  
441 those in non-disinfected systems. Third, non-disinfected systems exhibit greater diversity not only  
442 within a system (Figure 2) but also across systems as compared to disinfected systems. Despite the  
443 strong correlation between pairwise Mash distances of reads and scaffolds (Figure 3F), the median  
444 Mash distances for pairwise comparison of samples within each type of system (i.e., disinfected  
445 and non-disinfected) is higher for the scaffold-based analyses as compared to the read-based  
446 analyses. This is likely from the omission of low abundance microorganisms during *de novo*  
447 assembly and thus suggests that composition of medium-to-high abundance organisms are likely  
448 to be more variable between non-disinfected systems as compared to disinfected systems.

449 Finally, while the water chemistry and environmental parameters between disinfected and non-  
450 disinfected systems were distinct (Figure 1B), the parameters that most strongly correlated with  
451 Mash distances between samples were limited to a combination of chlorine, phosphate, and TOC  
452 for BioEnv analyses and chlorine, conductivity, and DO based on dbRDA. Both independent  
453 exploratory analyses consistently identified chlorine presence/absence and concentration as one of  
454 the key drivers of difference in microbial communities across the samples. Further, variance  
455 partition analyses indicated that ~17% of the variance in the Mash distance matrix was driven  
456 exclusively by chlorine; this make chlorine the most important parameter measured as part of this  
457 study in terms of differentiating between drinking water metagenomes. The significance of  
458 phosphate determined by BioEnv analyses is likely because chlorine and phosphate concentrations  
459 are inherently associated due to common use of the latter for corrosion control in DWDSs that  
460 maintain a chlorine residual<sup>68</sup>. Further, while it is unlikely that DO (identified as significant by  
461 dbRDA) directly affects microbial community composition (all DO concentrations were near or  
462 greater than saturation), it is possible that this may reflect the use of advanced oxidation process  
463 (e.g., ozonation) during drinking water treatment. Similarly, conductivity (identified as significant  
464 by dbRDA) is unlikely to directly influence the microbial community, but rather this may reflect  
465 the source water type and treatment processes being used for drinking water production.  
466 Specifically, source water derived from ground water sources or from reservoirs under the  
467 influence of ground water typically have much higher conductivities than those that rely on surface  
468 water supply. Similarly, chemicals used for softening and coagulation/flocculation processes may  
469 influence water conductivity. Thus, we speculate that the influence of conductivity may serve as

470 surrogate for a combination of source water and treatment process. These analyses clearly identify  
471 chlorine as one of the major measured parameters driving the Mash distances between samples,  
472 followed by conductivity (potential surrogate for source water and treatment process). Further, the  
473 fact the major proportion of the variance remains unexplained suggests that additional aspects such  
474 as treatment process configuration, DWDS characteristics, and other water chemistry parameters  
475 which were not characterized/measured as part of this study also likely play a strong role in  
476 differentiating between microbial communities in disinfected and non-disinfected drinking water  
477 systems.



**Figure 3:** Comparison of nucleotide composition using paired reads each from each sample and true scaffolds in each drinking water system according to Mash distance. (A, B) Heatmaps based on pairwise Mash distances of reads and scaffolds. Heatmaps are colored according to Mash distance; yellow denotes a distance of 0. Labels on x- and y-axis are colored according to disinfection strategies. (C). NMDS clustering of read based Mash distances between samples with vectors representing water chemistry/environmental parameters implemented using dbRDA. (D, E) Violin plots indicating the distribution of pairwise Mash distances of reads and scaffolds. Plots are colored according to the system type for which pairwise comparisons were conducted. Purple denotes comparisons between disinfected samples, yellow denotes comparisons between non-disinfected samples, and green denotes comparisons between disinfected and non-disinfected samples. (F) Correlation between average Mash distances of reads across samples and Mash distances of scaffolds across sampling locations.

478

479

480

481

482

483 **Protein coding sequences cluster based on presence/absence of disinfectant residuals.** A total  
484 of 8 million protein coding sequences were predicted and translated from true scaffolds, of which  
485 approximately 17 to 27% were annotated against KEGG database (Table S6). Consistent with the  
486 nucleotide-level analyses, samples clustered based on the presence and absence of disinfectant  
487 residual (Figure 4A, 4B, 4C) rather than by DWDS. Further, BioEnv analyses identified the  
488 combination of chlorine, phosphate, and ammonia as being strongly and significantly correlated  
489 ( $R = 0.392$ ,  $P < 0.001$ ) with Bray-Curtis distances between samples estimated using abundance (i.e.,  
490 RPKM) of KOs (Table S7). Similar to nucleotide based analyses, chlorine presence/absence and  
491 concentration was the measured parameter more strongly and significantly associated with  
492 differences in functional potential between samples at the single parameter level ( $R = 0.382$ ,  
493  $p < 0.001$ ). In contrast to nucleotide based analyses, conductivity and chlorine were the only two  
494 variables identified as significantly associated with Bray-Curtis distances between samples  
495 estimated using relative abundance of KO's in samples using dbRDA (Table S8). Variance  
496 partitioning indicated that both conductivity and chlorine individually explained approximately  
497 6.5% of the variance in Bray-Curtis distance matrix estimated using KO abundance. A comparison  
498 of the pairwise Mash distances within each group (i.e., disinfected, non-disinfected) and between  
499 them indicated that the diversity in functional potential was significantly different for both  
500 predicted protein coding-sequences and KEGG annotated proteins ( $p < 0.0001$ ). The median value  
501 of Mash distances between the non-disinfected samples was greater than that for disinfected  
502 samples (Figure 4D, 4E) and the differences in Mash distances between two groups was larger  
503 than the distances within each group. And finally, despite the fact that only 17-27% of predicted  
504 proteins were annotated against the KEGG database, the Mash distances between metagenomes  
505 estimated using all predicted protein coding sequences and those that were annotated against the  
506 KEGG database were highly correlated (Pearson's  $R \approx 1.00$ ,  $P < 0.05$ ) (Figure 4F), suggesting  
507 that focusing on annotated proteins does not result in significant loss of information while  
508 performing direct comparisons between samples from disinfected and non-disinfected systems.

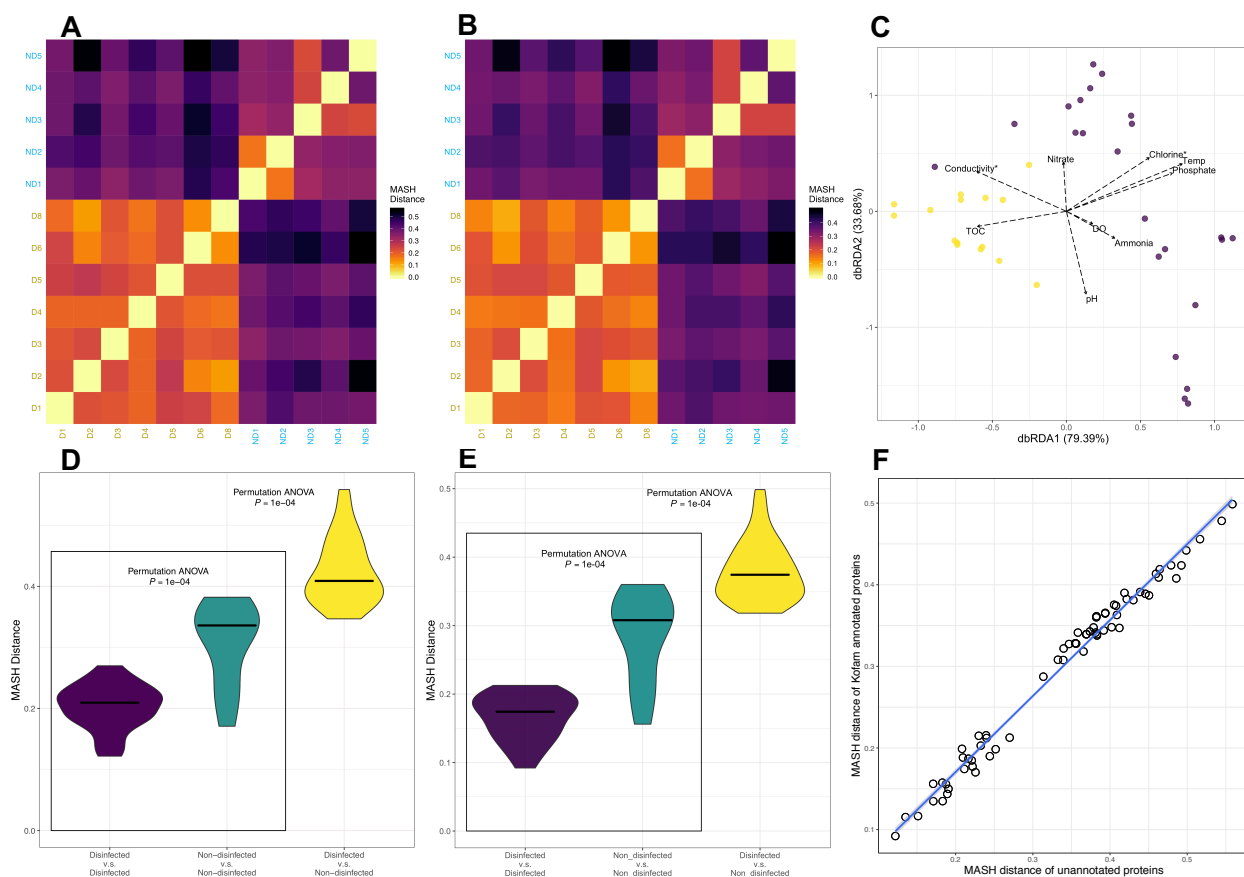
509 These analyses based on protein coding sequencing provide several key insights. First, clustering  
510 of samples into disinfected and non-disinfected groups is consistent for both community  
511 composition (i.e., read-based nucleotide composition analyses) and functional potential,  
512 irrespective of the use of all predicted ORF's and KEGG annotated protein sequences. Non-  
513 disinfected systems are significantly more heterogeneous across systems as compared to their

514 disinfected counterparts. This suggests that selection pressures exerted within disinfected systems  
515 are not only evident at community structure/membership (Figure 3), but also evident at the  
516 community functional potential level. Further, consistent with microbial community composition,  
517 chlorine was also identified as one of the key measured parameters driving differences between  
518 samples based on functional potential using both BioEnv and dbRDA analyses. In contrast to TOC  
519 which was included in the BioEnv parameter combination for microbial community composition  
520 level analyses, ammonia was identified as part of the combination at the functional potential level.  
521 While the exact reason behind this difference cannot be ascertained in this study, this may likely  
522 be associated with the fact that non-disinfected systems are severely nitrogen limited as compared  
523 to disinfected systems, while both systems were likely not carbon limited. Similar to the nucleotide  
524 level analyses, both conductivity and chlorine were identified as significantly ( $p < 0.01$ ) associated  
525 with differences between samples, with variance partitioning analyses allocating equal amount of  
526 variation to both parameters (Table S9). As speculated above, if conductivity is considered a signal  
527 for source water and treatment process type, then the impact of these two parameters on the  
528 functional potential of microbial community is relatively similar to that of presence/absence of the  
529 disinfectant residual. Finally, the residuals from the variance partitioning analyses were noticeably  
530 larger (84%) for functional potential analyses as compared to the microbial community  
531 composition (60%), suggesting that the impact of unmeasured/uncharacterized factors/parameters  
532 on microbial community functional potential was significantly larger than their impact on  
533 community composition. While it cannot be ruled out, it is unlikely that the higher fraction of  
534 unexplained variation was due to only a proportion of ORFs being annotated; this is because Mash  
535 distances estimated using only KEGG annotated ORFs were highly correlated with those estimated  
536 using all predicted ORFs using suggesting little to minimal loss of discriminatory power while  
537 using only annotated proteins.

538

539

540



**Figure 4:** Comparison of functional potential among all and KEGG protein-coding amino acid sequences across sampling locations. This analysis estimates dissimilarity in amino acid composition of samples, similar to the nucleotide composition analyses presented earlier. (A, B) Heatmaps based on pairwise Mash distances of all protein coding sequences and Bray-Curtis distances using KO. Heatmaps are colored according to Mash/Bray-Curtis distance; yellow denotes a distance of 0. Labels on x- and y-axis are colored according to disinfection strategies; dark golden denotes samples with chlorine, while blue denotes samples without disinfectant residuals. C). NMDS clustering of using Bray-Curtis distances using KO abundances between samples with vectors representing water chemistry/environmental parameters implemented using dbRDA. Violin plots indicating the distribution of pairwise (D) Mash distances of all and (E) Bray-Curtis distances KEGG annotated proteins. Crossbars indicate the median value of Mash distances. (F) Correlation between pairwise Mash distances estimated using all and Bray-Curtis distances for KEGG annotated protein coding sequences.

541

542 **Differentially abundant metabolic modules are consistent with microbial growth control**  
 543 **strategies.** A total of 7,281 KOs were identified in all samples with 5,922 remaining post-filtering  
 544 based on scaffold coverage (>1x) and frequency of KO detection in each drinking water system  
 545 (detected more than once) (Table S10). The 5,922 KO's were further categorized into 540 KEGG  
 546 modules and upon further filtering to remove KEGG modules with no more than one missing block  
 547 and greater than equal to 50% completion, a total of 208 KEGG modules were retained (Table



548 S11). Of these, a total of 57 KEGG modules exhibited significantly differential abundance between  
549 disinfected and non-disinfected samples ( $p$ -value < 0.005) (Table S12, S13). Modules associated  
550 with ribosomal synthesis, ribonucleotide biosynthesis, and RNA polymerase were ignored from  
551 further consideration. Similarly, modules most likely associated with plant metabolism (e.g.,  
552 Crassulacean acid metabolism) were also ignored. This resulted in 29 and 22 KEGG modules that  
553 were more abundant in non-disinfected system and disinfected systems, respectively. These  
554 included modules associated with energy metabolism (disinfected, i.e. D=2, non-disinfected, i.e.,  
555 ND=5), carbohydrate and lipid metabolism (D=11, ND=10), nucleotide and amino acid  
556 metabolism (D=5, ND=13), and secondary metabolism (D=4, ND=1).

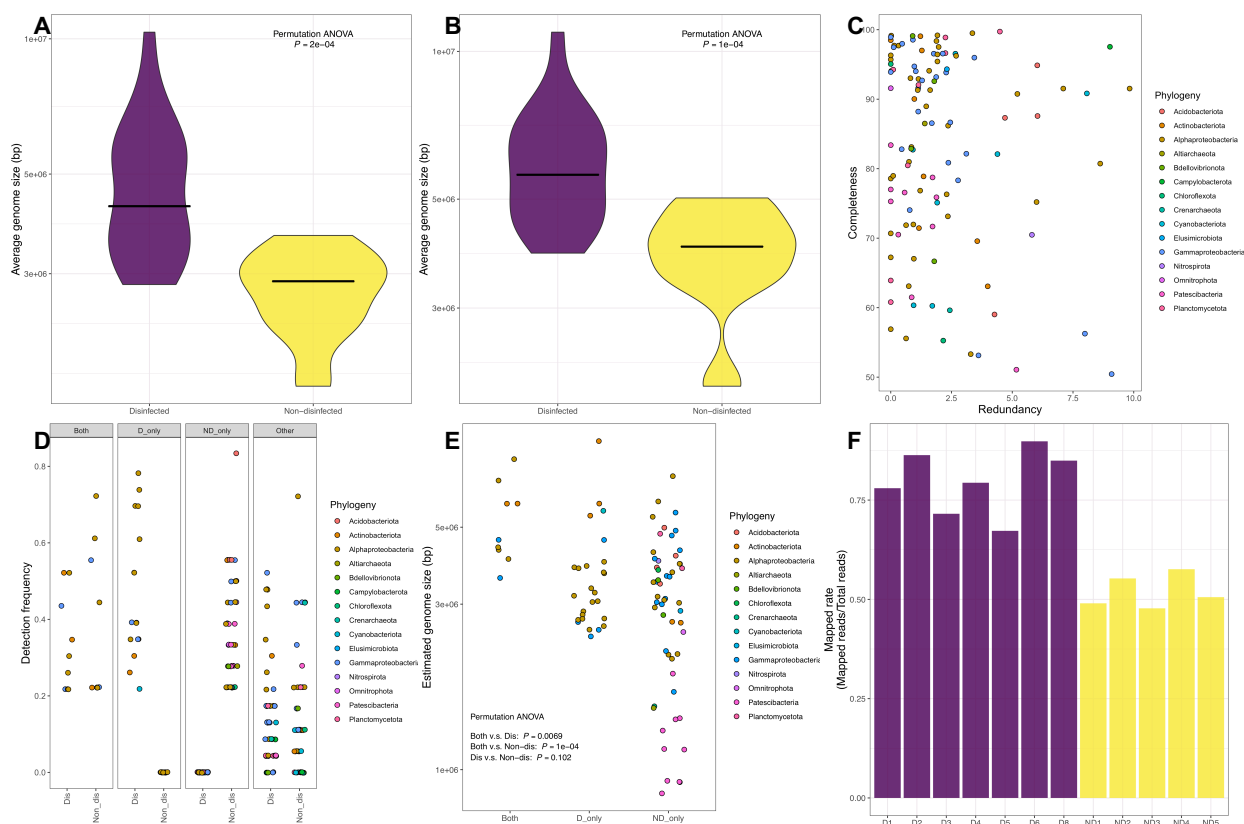
557 Metabolic modules associated with polyamine biosynthesis, aromatics degradation, terpenoid  
558 biosynthesis, and fatty acid metabolism were significantly enriched in disinfected systems.  
559 Specifically, metabolic pathways associated with benzene (M00548) and benzoate (M00551)  
560 degradation to catechol and methyl catechol were highly enriched in disinfected systems. Further,  
561 eukaryota-associated metabolic modules such as terpenoid backbone biosynthesis (M00367) and  
562 modules associated with peroxisomal beta-oxidation of very long chain fatty acids (M00861) are  
563 likely to be enriched in the disinfected systems due to the higher relative abundance of eukaryota  
564 in samples collected from disinfected as compared to non-disinfected systems respectively.  
565 Further, modules related to  $\gamma$ -aminobutyrate (GABA) metabolism (M00136, M00027) were  
566 enriched in disinfected systems. The GABA shunt pathway converts glutamate to GABA using  
567 glutamate decarboxylase (GAD), followed by reversible conversion from  $\alpha$ -ketoglutarate to  
568 succinate semialdehyde (SSA) through the activity of GABA transaminase (GABA-AT), and  
569 finally succinate is formed by succinate semialdehyde dehydrogenase (SSDH) activity. In contrast,  
570 the key metabolic modules enriched in non-disinfected systems were associated with carbon  
571 fixation and methane metabolism (M00377, M00620, and M00422) and nitrogen fixation  
572 (M00175) (Table S13). The differentially abundant carbon fixation modules included the Wood-  
573 Ljungdahl pathway, Acetyl-CoA pathway, and the incomplete reductive citrate cycle. These  
574 pathways can fix carbon dioxide to produce acetyl-CoA which can then be converted to other  
575 necessary biosynthetic intermediates of the carbon metabolism<sup>70, 71</sup>.

576 The enrichment of carbon and nitrogen fixation modules in non-disinfected systems is consistent  
577 with nutrient limitation as the strategy for microbial growth control in non-disinfected drinking

578 water systems. While the measured total organic carbon concentrations in non-disinfected systems  
579 did not indicate carbon limited conditions, DWTP's supplying water to non-disinfected DWDSs  
580 typically achieve far superior levels of removal of assimilable organic carbon (AOC)<sup>28</sup>. Similarly,  
581 the nitrogen availability in the form of ammonia was consistently zero for non-disinfected systems  
582 compared to disinfected systems which has residual ammonia concentrations ranged from 0.01-  
583 0.15 mg/l of ammonia-nitrogen. In contrast, the enrichment of KEGG modules associated with  
584 GABA metabolism in disinfected systems suggests the potential importance of stress protection  
585 and utilization of microbial decay products. Previous studies have shown that GABA metabolism  
586 is associated with bacterial survival under various types of environmental stresses, including  
587 oxidative stress, acidic stress, and osmotic stress<sup>72-75</sup>. Meanwhile, GABA can also play a  
588 significant role in nitrogen metabolism of bacteria. For instance, putrescine formed due to the  
589 breakdown of amino acids potentially from decaying biomass, can be converted to GABA  
590 (M00136) and finally metabolized via GABA shunt pathway<sup>74</sup>. The enrichment of GABA  
591 metabolism in disinfected systems may thus be associated with greater protection against  
592 disinfectant stress and by allowing access to decay products from inactivated cells.

593 **Average genome size differences between disinfected and non-disinfected system vary**  
594 **between read-based and MAG-based analyses.** We further investigated differences in genome  
595 sizes between disinfected and non-disinfected systems. Genome sizes can be indicative of the  
596 metabolic capacity of microorganisms<sup>76</sup> and thus provide insights in the whether the  
597 presence/absence of disinfectants selects for organisms with larger or smaller metabolic  
598 repertoire<sup>77</sup> in comparison to organisms detected in non-disinfected systems. Average genomes  
599 size estimates from disinfected systems were significantly larger than those from non-disinfected  
600 systems based on MicrobeCensus estimates using entire metagenomic data (Figure 5A); this was  
601 consistent even when reads mapping to phyla known to have smaller genomes (e.g.,  
602 *Patescibacteria*) were selectively removed from the data set (Figure 5B). This suggests that  
603 microorganisms in disinfected systems may be metabolically more diverse than their counterparts  
604 from non-disinfected systems. Nonetheless, these results were not consistent when compared with  
605 estimated genome sizes of MAGs recovered as part of this study. Specifically, we recovered a total  
606 of 115 dereplicated MAGs with completeness >50% and redundancy <10% (Table S14). These  
607 115 MAGs were binned into four categories based on the detection or non-detection in disinfected  
608 samples. Specifically, MAGs were binned in the four groups (i.e., both, D-only, ND-only, and

609 other) based on genome coverage and detection frequency criteria outlined in the materials and  
610 methods section (see MAG-level analyses) (Table S15). This resulted in 9, 16, 41, and 49 MAGs  
611 were categorized as both, D-only, ND-only, and other (Figure 5C, 5D) (Table S14). In contrast to  
612 read-based estimates of average genome size, MAG-based genome size estimates were not  
613 significantly different between the three key categories (Both= $4.4\pm 0.77$ Mbp, D-  
614 only= $3.22\pm 0.81$ Mbp, ND-only= $3.48\pm 1.22$ Mbp) (Figure 5E). Yet, the ND-only category  
615 consisted of several smaller genomes (n=17) compared to the D category. The lack of genome size  
616 differences between disinfected and non-disinfected samples based on MAG-based analyses  
617 compared to metagenome-level read-based analyses may be due to the proportion of read-based  
618 data represented by the MAGs. Specifically, while 60-90% of the reads from disinfected systems  
619 mapped to the 115 MAGs with the mapping rate from non-disinfected systems averaging around  
620 50% (Figure 5F). Thus, it is likely that the metagenomic assembly and binning process may have  
621 resulted in suboptimal recovery of smaller genomes from non-disinfected sample which eliminates  
622 the signal in genome size differences observed at the metagenome level.



**Figure 5:** Violin plots indicating the genome size estimated by MicrobeCensus (a) before and (b) after Patescibacteria removal suggest average genome sizes in disinfected systems are larger than those in non-disinfected systems. (C) The 115 MAGs assembled with >50% completeness and <10% redundancy were categorized into (D) four groups based on their detection frequency in disinfected and non-disinfected systems. (E) While the estimated genome sizes of MAGs in D\_only, ND\_only, and Both categories were not significantly different, the ND\_only category consisted of large number of smaller genomes. (F) Barplot indicating the proportion of reads mapped to 115 genomes across samples. Purple and yellow denotes samples from systems with and without a disinfectant residual, respectively.

623 **Metabolic capacities differ between metagenome assembled genomes from disinfected and**  
 624 **non-disinfected systems.** Clustering of MAGs (Figure 6A) based on presence/absence of KEGG  
 625 metabolic modules was largely driven by phylogenetic placement of MAGs, rather than their  
 626 classifications into groups based on the detection frequencies in disinfected and non-disinfected  
 627 systems (Figure 6B). Further, there was insufficient representation of MAGs from D-only/ND-  
 628 only categories across all phylogenetic clusters (e.g., at the species or genus level) to allow for  
 629 direct comparisons of metabolic potential of closely related MAGs exclusively frequent in  
 630 disinfected and non-disinfected systems. Nonetheless, there were seven and five high quality  
 631 (completeness > 90%, redundancy <10%) alphaproteobacterial MAGs that were exclusively  
 632 frequent in disinfected (average detection frequency in disinfected =55%) and non-disinfected  
 633 systems (average detection frequency in non-disinfected=29%) (Figure 6A). Thus, we focused

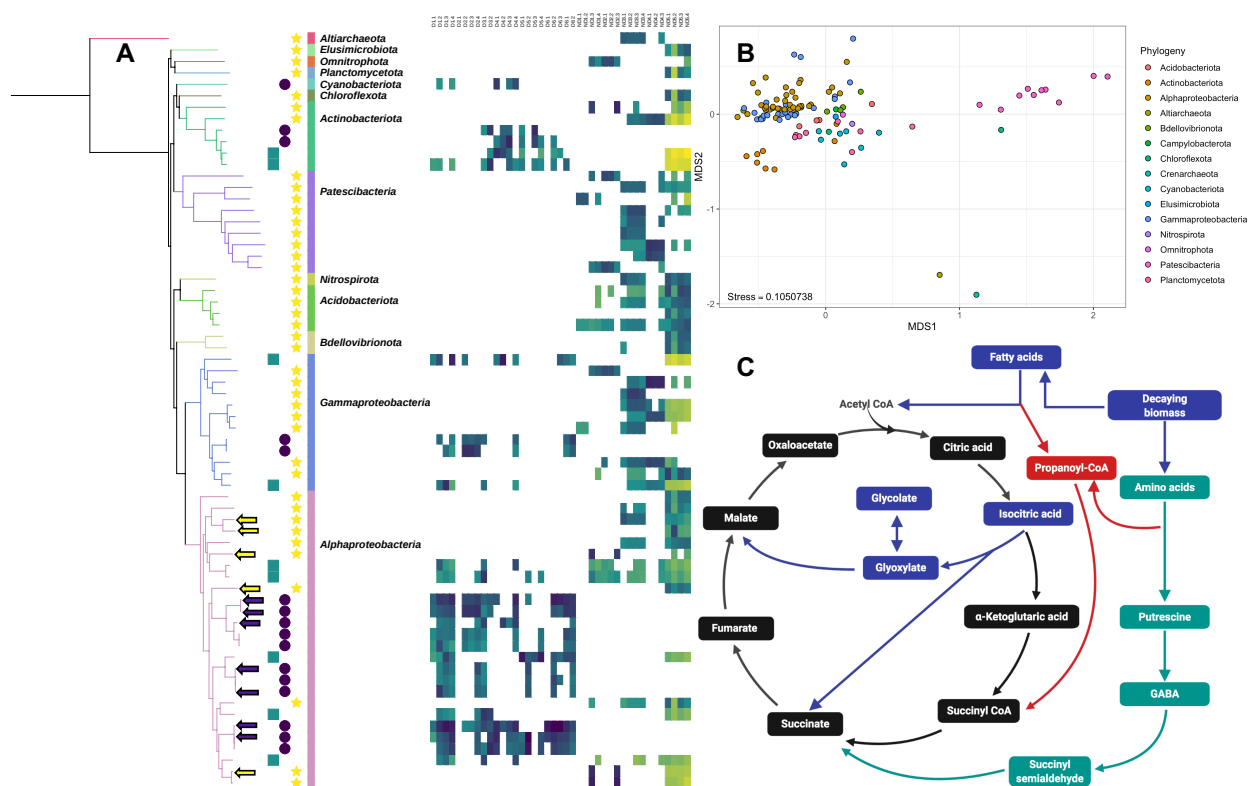
634 metabolic module comparisons between these 12 MAGs only. We evaluated differences in  
635 metabolic capacity of these MAGs by (1) considering all KEGG modules  $\geq 75\%$  complete within  
636 MAGs to be present in them and (2) all modules present in more than half of the high-quality  
637 MAGs within each category to be present within each category (Figure 6C, Table S16). We  
638 subsequently confirmed the presence/absence of genes within key metabolic modules using KO-  
639 level annotation for these 12 MAGs (Table S17).

640 The metabolic module associated with the glyoxylate cycle (M00012) was present in 86% of the  
641 MAGs in the D-only category while being only partially complete in most of the ND-only MAGs.  
642 Specifically, isocitrate lyase (*aceA*: K01637) and malate synthase (*aceB*: K01638), two key genes  
643 involved in the glyoxylate cycle, were present in 40% and 100% of the MAGs from D-only,  
644 respectively and both genes were absent in all ND-only MAGs included in this analysis. The  
645 glyoxylate shunt is associated with use of non-carbohydrate carbon sources (i.e., via  
646 gluconeogenesis), such as break down products from lipids, fatty acids etc<sup>78</sup>. The likely benefit of  
647 the glyoxylate shunt and associated use of lipids and fatty acids as carbon source is further  
648 supported by the fact that KEGG module associated with propanoyl-coA metabolism (M00741)  
649 was complete in 6/7 as compared to 2/5 MAGs from the D-only and ND-only categories. This  
650 metabolic module is associated with the conversion of propanoyl-coA, a toxic byproduct of fatty  
651 and amino acid degradation, to succinyl-coA. High biomass turnover rates, due to disinfectant  
652 induced microbial inactivation, may result in resource pools enriched in microbial decay products  
653 thus allowing a significant advantage for microorganisms capable of necrotrophic growth<sup>79</sup> aided  
654 by the glyoxylate cycle. Thus, it is feasible that the ability to utilize microbial decay products may  
655 provide a distinct advantage to microorganisms inhabiting disinfected DWDSs.

656 The glyoxylate shunt may provide additional benefits for microorganisms subject to disinfectant  
657 stress via enhanced fitness to oxidative stress<sup>78</sup> and enhanced persistence when challenged with  
658 other chemical stressors (e.g., antibiotics)<sup>80</sup>. In contrast to module level analyses at the  
659 metagenome level where carbon fixation capacity was significantly more abundant in non-  
660 disinfected as compared to disinfected systems, the alphaproteobacterial MAGs from D-only  
661 systems harbored the capacity for carbon fixation via the Calvin-Benson-Bassham cycle (M00165,  
662 M00166, M00167) while this capacity was mostly absent from MAGs in the ND-only category.  
663 Nonetheless, these MAG-based analyses are limited in phylogenetic scope and does not weigh the

664 importance of MAGs to their respective systems based on their relative abundance. Hence, we  
665 suggest that metagenome-level analyses should take precedence over findings at the MAG level  
666 when they conflict. While the glyoxylate shunt was not identified as a significantly enriched in the  
667 disinfected systems at the metagenome level analyses, the GABA shunt (metagenome level  
668 analyses) and glyoxylate shunt (MAG level analyses) may both be involved in use of non-  
669 carbohydrate carbon sources suggesting that re-use of microbial decay products may indeed be a  
670 key bacterial trait that allows for persistence in disinfected drinking water systems. Further lending  
671 support to this is that that propanoyl-coA metabolism was identified as significantly enriched in  
672 disinfected systems compared to non-disinfected systems using both metagenome-level and MAG-  
673 level analyses. Interestingly, only one metabolic module was identified as being more than twice  
674 as prevalent in alphaproteobacterial MAGs from ND-only systems compared to those from D-only  
675 systems (i.e., M00156: cbb3-type Cytochrome C oxidase). The greater metabolic capacity of  
676 alphaproteobacterial D-only MAGs compared to ND-only MAGs was also confirmed at the KO-  
677 level by evaluating the presence/absence of KO's in the D-only and ND-only category MAGs.  
678 Specifically, while only 8 KOs were twice or more as prevalent in ND-only MAGs compared to  
679 D-only MAGs, the total KOs that were twice or more as prevalent in D-only MAGs was 109. This  
680 supports the conclusion that metabolic repertoire of alphaproteobacterial D-only MAGs is  
681 significantly larger than that of ND-only MAGs. Notable among the genes that were twice as  
682 frequent in D-only MAGs compared ND-only MAGs included those involve in SOS-response  
683 mediated mutagenesis involving trans-lesion synthesis (i.e., *imuA*: K14160, *imuB*: K14161, and  
684 *dnaE2*: K14162)<sup>81</sup>, glyoxylate reductase (*gyaR*: K00015) which may be likely involved in  
685 regulating glyoxylate concentrations, and vitamin B12 transporter (*btuB*: K16092). SOS response  
686 is typically activated in response to significant cellular accumulation of damaged DNA<sup>82</sup> and *imuA*  
687 and *imuB* co-expression with *dnaE2* has been shown to be responsive to UV damage<sup>81</sup>. Thus, the  
688 higher prevalence of SOS response related genes in D\_only MAGs may be associated with the  
689 DNA damage caused by disinfectants. Further, the ability to synthesize vitamin B12, an essential  
690 co-factor, is limited to certain bacteria and archaea and thus the ability to uptake vitamin B12 from  
691 the environment is essential for growth<sup>83</sup>. The higher abundance of vitamin B12 transporters is  
692 consistent with metagenome level observations that the microbial community in disinfected  
693 systems rely more on scavenging from the environment as compared to non-disinfected systems.

694



**Figure 6:** (A) Phylogenomic tree of 66 MAGs classified as D-only (purple circles), ND-only (yellow stars), and both (teal squares) constructed using 48 ribosomal proteins along and their relative abundance (RPKM) in the samples collected from disinfected and non-disinfected systems. RPKM's for MAGs are only reported for samples where 25% of the nucleotides in a MAG were covered by at least one read. (B) Clustering of all MAGs based on their clustering metabolic potential (i.e., completeness of KEGG modules) was primarily drive by phylogeny. (C) The metabolic modules identified as differentially abundant in disinfected systems using metagenome level analyses (Table S12) are shown using teal arrows and squares and those more prevalent in high quality alphaproteobacterial MAGs from D-only (purple arrows - Figure 6A) compared to those from ND-only category (yellow arrows - Figure 6A) are shown using blue arrows and boxes, while red arrows and boxes indicates modules identified as more prevalent in D-only systems using both metagenome and MAG level analyses.

695

## 696 CONCLUSIONS.

697 To our knowledge, this is the first study to provide metagenomic insights into differences in  
698 structure and functional potential of drinking water microbiomes across full-scale drinking water  
699 systems that rely on disinfection (i.e., disinfected) or nutrient limitation (i.e., non-disinfected) to  
700 manage microbial growth. Understanding the microbial implications of these two microbial  
701 growth control strategies is essential to not only develop a better understanding of ecological and  
702 metabolic traits guiding community level processes in these system, but is also critical for  
703 providing a community-level context to the microbiological safety in either type of drinking water  
704 system. In this study, we show that disinfection exhibits consistent, systematic, and significant  
705 association with drinking water microbiome at the membership, structure, and functional potential  
706 at the metagenome and MAG levels, irrespective of the drinking water system under consideration  
707 (e.g., source water type, treatment process, etc.). In doing so, we also identify key metabolic traits  
708 associated with carbon and nitrogen metabolism that are over represented in bacteria in disinfected  
709 systems compared to non-disinfected systems. This suggests that the influence and efficacy of  
710 disinfection on the drinking water microbiome may not simply be associated with differential  
711 disinfection resistance<sup>84</sup>, but may also expand to other metabolic traits that include the use of  
712 carbon and nitrogen sources made available via microbial inactivation and its regulation. It is  
713 important to note that while the impact of disinfection on microbial community structure and  
714 functional potential is clear, the metabolic traits identified in this study provide a hypothesis to  
715 support future experimental work that will be required to validate the findings of this study.

## 716 ACKNOWLEDGEMENTS

717 ZD was supported by the Lord Kelvin Adam Smith Scholarship at the University of Glasgow.  
718 MS was supported by the College of Engineering at Northeastern University. This work was  
719 funded by Engineering and Physical Science Research Council (EP/M016811/1) and the  
720 National Science Foundation (NSF-CBET 1749530). UZI is supported by NERC Independent  
721 Research Fellowship (NERC NE/L011956/1). The authors are grateful to Prof. Karthik  
722 Anantharaman for providing helpful critiques of the manuscript.

723

724



725 REFERENCES

726

- 727 1. LeChevallier, M. W.; Welch, N. J.; Smith, D. B., Full-scale studies of factors related to  
728 coliform regrowth in drinking water. *Applied and Environmental Microbiology* **1996**, *62*,  
729 (7), 2201-2211.
- 730 2. Liu, G.; Bakker, G. L.; Li, S.; Vreeburg, J. H. G.; Verberk, J. Q. J. C.; Medema, G. J.; Liu,  
731 W. T.; Van Dijk, J. C., Pyrosequencing Reveals Bacterial Communities in Unchlorinated  
732 Drinking Water Distribution System: An Integral Study of Bulk Water, Suspended Solids,  
733 Loose Deposits, and Pipe Wall Biofilm. *Environmental Science & Technology* **2014**, *48*,  
734 5467-5476.
- 735 3. Liu, G.; Zhang, Y.; van der Mark, E.; Magic-Knezev, A.; Pinto, A.; van den Bogert, B.; Liu,  
736 W.; van der Meer, W.; Medema, G., Assessing the origin of bacteria in tap water and  
737 distribution system in an unchlorinated drinking water system by SourceTracker using  
738 microbial community fingerprints. *Water Research* **2018**, *138*, 86-96.
- 739 4. Berry, D.; Xi, C.; Raskin, L., Microbial ecology of drinking water distribution systems.  
740 *Current Opinion in Biotechnology* **2006**, *17*, (3), 297-302.
- 741 5. Proctor, C. R.; Hammes, F., Drinking water microbiology—from measurement to  
742 management. *Current Opinion in Biotechnology* **2015**, *33*, 87-94.
- 743 6. Chao, Y.; Ma, L.; Yang, Y.; Ju, F.; Zhang, X.-X.; Wu, W.-M.; Zhang, T., Metagenomic  
744 analysis reveals significant changes of microbial compositions and protective functions  
745 during drinking water treatment. *Scientific Reports* **2013**, *3*, 3550.
- 746 7. Roeselers, G.; Coolen, J.; van der Wielen, P. W. J. J.; Jaspers, M. C.; Atsma, A.; de Graaf,  
747 B.; Schuren, F., Microbial biogeography of drinking water: Patterns in phylogenetic  
748 diversity across space and time. *Environ. Microbiol.* **2015**, *17*, (7), 2505–2514.
- 749 8. Pinto, A. J.; Xi, C.; Raskin, L., Bacterial Community Structure in the Drinking Water  
750 Microbiome Is Governed by Filtration Processes. *Environmental Science & Technology*  
751 **2012**, *46*, (16), 8851-8859.
- 752 9. Lautenschlager, K.; Hwang, C.; Ling, F.; Liu, W. T.; Boon, N.; Köster, O.; Egli, T.;  
753 Hammes, F., Abundance and composition of indigenous bacterial communities in a multi-  
754 step biofiltration-based drinking water treatment plant. *Water Res.* **2014**, *62*, 40–52.
- 755 10. Pinto, A. J.; Schroeder, J.; Lunn, M.; Sloan, W.; Raskin, L., Spatial-Temporal Survey and  
756 Occupancy-Abundance Modeling To Predict Bacterial Community Dynamics in the  
757 Drinking Water Microbiome. *mBio* **2014**, *5*, (3), e01135-14.
- 758 11. Baron, J. L.; Vikram, A.; Duda, S.; Stout, J. E.; Bibby, K., Shift in the Microbial Ecology of  
759 a Hospital Hot Water System following the Introduction of an On-Site Monochloramine  
760 Disinfection System. *PLoS ONE* **2014**, *9*, (7).
- 761 12. Proctor, C. R.; Gächter, M.; Kötzsch, S.; Rölli, F.; Sigrist, R.; Walser, J.-C.; Hammes, F.,  
762 Biofilms in shower hoses – choice of pipe material influences bacterial growth and  
763 communities. *Environmental Science: Water Research & Technology* **2016**, *2*, (4), 670-682.
- 764 13. Jia, S.; Shi, P.; Hu, Q.; Li, B.; Zhang, T.; Zhang, X. X., Bacterial Community Shift Drives  
765 Antibiotic Resistance Promotion during Drinking Water Chlorination. *Environ. Sci.*  
766 *Technol.* **2015**, *49*, (20), 12271–12279.
- 767 14. Wang, H.; Masters, S.; Edwards, M. A.; Falkinham, J. O.; Pruden, A., Effect of  
768 Disinfectant, Water Age, and Pipe Materials on Bacterial and Eukaryotic Community  
769 Structure in Drinking Water Biofilm. *Environmental Science & Technology* **2014**, *48*, (3),  
770 1426-1435.

- 771 15. Prest, E. I.; Hammes, F.; van Loosdrecht, M. C. M.; Vrouwenvelder, J. S., Biological  
772 Stability of Drinking Water: Controlling Factors, Methods, and Challenges. *Frontiers in*  
773 *Microbiology* **2016**, *7*.
- 774 16. Wang, H.; Pryor, M. A.; Edwards, M. A.; Falkinham, J. O.; Pruden, A., Effect of GAC pre-  
775 treatment and disinfectant on microbial community structure and opportunistic pathogen  
776 occurrence. *Water Research* **2013**, *47*, (15), 5760-5772.
- 777 17. Liu, G.; Verberk, J. Q. J. C.; Van Dijk, J. C., Bacteriology of drinking water distribution  
778 systems: an integral and multidimensional review. *Applied Microbiology and Biotechnology*  
779 **2013**, *97*, (21), 9265-9276.
- 780 18. van der Wielen, P. W. J. J.; van der Kooij, D., Nontuberculous Mycobacteria, Fungi, and  
781 Opportunistic Pathogens in Unchlorinated Drinking Water in the Netherlands. *Applied and*  
782 *Environmental Microbiology* **2013**, *79*, (3), 825.
- 783 19. Garner, E.; McLain, J.; Bowers, J.; Engelthaler, D. M.; Edwards, M. A.; Pruden, A.,  
784 Microbial Ecology and Water Chemistry Impact Regrowth of Opportunistic Pathogens in  
785 Full-Scale Reclaimed Water Distribution Systems. *Environ Sci Technol* **2018**, *52*, (16),  
786 9056-9068.
- 787 20. van Lieverloo, J. H. M.; Hoogenboezem, W.; Veenendaal, G.; van der Kooij, D., Variability  
788 of invertebrate abundance in drinking water distribution systems in the Netherlands in  
789 relation to biostability and sediment volumes. *Water Research* **2012**, *46*, (16), 4918-4932.
- 790 21. Christensen, S. C. B.; Nissen, E.; Arvin, E.; Albrechtsen, H.-J., Distribution of *Asellus*  
791 *aquaticus* and microinvertebrates in a non-chlorinated drinking water supply system –  
792 Effects of pipe material and sedimentation. *Water Research* **2011**, *45*, (10), 3215-3224.
- 793 22. Otten, T. G.; Graham, J. L.; Harris, T. D.; Dreher, T. W., Elucidation of Taste- and Odor-  
794 Producing Bacteria and Toxigenic Cyanobacteria in a Midwestern Drinking Water Supply  
795 Reservoir by Shotgun Metagenomic Analysis. *Appl. Environ. Microbiol.* **2016**, *82*, (17),  
796 5410-5420.
- 797 23. Beech, I. B.; Sunner, J., Biocorrosion: towards understanding interactions between biofilms  
798 and metals. *Current Opinion in Biotechnology* **2004**, *15*, (3), 181-186.
- 799 24. Zhang, Y.; Griffin, A.; Edwards, M., Nitrification in Premise Plumbing: Role of Phosphate,  
800 pH and Pipe Corrosion. *Environmental Science & Technology* **2008**, *42*, (12), 4280-4284.
- 801 25. Rosario-Ortiz, F.; Rose, J.; Speight, V.; Gunten, U. v.; Schnoor, J., How do you like your  
802 tap water? *Science* **2016**, *351*, (6276), 912–914.
- 803 26. Potgieter, S.; Pinto, A.; Sigudu, M.; du Preez, H.; Ncube, E.; Venter, S., Long-term spatial  
804 and temporal microbial community dynamics in a large-scale drinking water distribution  
805 system with multiple disinfectant regimes. *Water Research* **2018**, *139*, 406-419.
- 806 27. Kooij, D. v. d.; Wielen, P. W. J. J. v. d.; Rosso, D.; Shaw, A.; Borchardt, D.; Ibsch, R.;  
807 Apgar, D.; Witherspoon, J.; Toro, D. M. d.; Paquin, P. R.; Mavinic, D.; Koch, F.; Guillot,  
808 E.; Loret, J.-F.; Hoffmann, E.; Ødegaard, H.; Hernandez-Sancho, F.; Molinos-Senante, M.,  
809 *Microbial Growth in Drinking Water Supplies*. IWA Publishing: 2013; p 484.
- 810 28. van der Kooij, D.; van der Wielen, P. W. J. J., Microbial Growth in Drinking-Water  
811 Supplies: Problems, Causes, Control and Research Needs. In IWA Publishing: 2013.
- 812 29. Richardson, S. D., Disinfection by-products and other emerging contaminants in drinking  
813 water. *TrAC Trends in Analytical Chemistry* **2003**, *22*, (10), 666-684.
- 814 30. Sedlak, D. L.; von Gunten, U., The Chlorine Dilemma. *Science* **2011**, *331*, (6013), 42–43.

- 815 31. Li, X.-F.; Mitch, W. A., Drinking Water Disinfection Byproducts (DBPs) and Human  
816 Health Effects: Multidisciplinary Challenges and Opportunities. *Environmental Science &*  
817 *Technology* **2018**, *52*, (4), 1681-1689.
- 818 32. Falkinham, J. O.; Pruden, A.; Edwards, M., Opportunistic Premise Plumbing Pathogens:  
819 Increasingly Important Pathogens in Drinking Water. *Pathogens* **2015**, *4*, (2), 373-386.
- 820 33. Zhang, H.; Chang, F.; Shi, P.; Ye, L.; Zhou, Q.; Pan, Y.; Li, A., Antibiotic Resistome  
821 Alteration by Different Disinfection Strategies in a Full-Scale Drinking Water Treatment  
822 Plant Deciphered by Metagenomic Assembly. *Environmental Science & Technology* **2019**,  
823 *53*, (4), 2141-2150.
- 824 34. Shi, P.; Jia, S.; Zhang, X. X.; Zhang, T.; Cheng, S.; Li, A., Metagenomic insights into  
825 chlorination effects on microbial antibiotic resistance in drinking water. *Water Research*.  
826 **2013**, *47*, (1), 111-120.
- 827 35. Sevillano, M.; Dai, Z.; Calus, S.; Santos, Q. M. B.-d. I.; Eren, A. M.; Wielen, P. W. J. J. v.  
828 d.; Ijaz, U. Z.; Pinto, A. J., Disinfectant residuals in drinking water systems select for  
829 mycobacterial populations with intrinsic antimicrobial resistance. *bioRxiv* **2019**, 675561.
- 830 36. Bertelli, C.; Courtois, S.; Rosikiewicz, M.; Piriou, P.; Aeby, S.; Robert, S.; Loret, J.-F.;  
831 Greub, G., Reduced Chlorine in Drinking Water Distribution Systems Impacts Bacterial  
832 Biodiversity in Biofilms. *Frontiers in Microbiology* **2018**, *9*.
- 833 37. Hamsch, B.; Böckle, K.; van Lieverloo, J. H. M., Incidence of faecal contaminations in  
834 chlorinated and non-chlorinated distribution systems of neighbouring European countries.  
835 *Journal of Water and Health* **2007**, *5*, (S1), 119-130.
- 836 38. Bautista-de los Santos, Q. M.; Schroeder, J. L.; Blakemore, O.; Moses, J.; Haffey, M.;  
837 Sloan, W.; Pinto, A. J., The impact of sampling, PCR, and sequencing replication on  
838 discerning changes in drinking water bacterial community over diurnal time-scales. *Water*  
839 *Research* **2016**, *90*, 216-224.
- 840 39. Waak, M. B.; Hozalski, R. M.; Hallé, C.; LaPara, T. M., Comparison of the microbiomes of  
841 two drinking water distribution systems with and without residual chloramine disinfection.  
842 *Microbiome* **2019**, *7*, (1), 87.
- 843 40. Clesceri, L. S.; Greenberg, A. E.; Eaton, A. D., *Standard Methods for the Examination of*  
844 *Water and Wastewater, 20th Edition*. APHA American Public Health Association: 1998.
- 845 41. Bolger, A. M.; Lohse, M.; Usadel, B., Trimmomatic: a flexible trimmer for Illumina  
846 sequence data. *Bioinformatics* **2014**, *30*, (15), 2114-2120.
- 847 42. Rodriguez-R, L. M.; Gunturu, S.; Tiedje, J. M.; Cole, J. R.; Konstantinidis, K. T., Nonpareil  
848 3: Fast Estimation of Metagenomic Coverage and Sequence Diversity. *mSystems* **2018**, *3*,  
849 (3).
- 850 43. Nayfach, S.; Pollard, K. S., Average genome size estimation improves comparative  
851 metagenomics and sheds light on the functional ecology of the human microbiome. *Genome*  
852 *Biology* **2015**, *16*, (1), 51.
- 853 44. Parks, D. H.; Chuvochina, M.; Waite, D. W.; Rinke, C.; Skarshewski, A.; Chaumeil, P. A.;  
854 Hugenholtz, P., A standardized bacterial taxonomy based on genome phylogeny  
855 substantially revises the tree of life. *Nature Biotechnol.* **2018**, *36*, (10), 996.
- 856 45. Nurk, S.; Meleshko, D.; Korobeynikov, A.; Pevzner, P. A., MetaSPAdes: A new versatile  
857 metagenomic assembler. *Genome Res.* **2017**, *27*, (5), 824-834.
- 858 46. Li, H., Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.  
859 *arXiv:1303.3997 [q-bio]* **2013**.

- 860 47. Ondov, B. D.; Treangen, T. J.; Melsted, P.; Mallonee, A. B.; Bergman, N. H.; Koren, S.;  
861 Phillippy, A. M., Mash: fast genome and metagenome distance estimation using MinHash.  
862 *Genome Bioinform.* **2016**, *17*, (1), 132.
- 863 48. Hyatt, D.; Chen, G. L.; LoCascio, P. F.; Land, M. L.; Larimer, F. W.; Hauser, L. J.,  
864 Prodigal: Prokaryotic gene recognition and translation initiation site identification. *BMC*  
865 *Bioinformatics* **2010**, *11*, 119.
- 866 49. Eddy, S. R., Accelerated profile HMM searches. *PLoS Computational Biology* **2011**, *7*, (10),  
867 e1002195.
- 868 50. El-Gebali, S.; Mistry, J.; Bateman, A.; Eddy, S. R.; Luciani, A.; Potter, S. C.; Qureshi, M.;  
869 Richardson, L. J.; Salazar, G. A.; Smart, A.; Sonnhammer, E. L. L.; Hirsh, L.; Paladin, L.;  
870 Piovesan, D.; Tosatto, S. C. E.; Finn, R. D., The Pfam protein families database in 2019.  
871 *Nucleic Acids Research*. **2019**, *47*, (D1), D427–D432.
- 872 51. Nawrocki, E. P.; Eddy, S. R., Infernal 1.1: 100-fold faster RNA homology searches.  
873 *Bioinformatics* **2013**, *29*, (22), 2933–2935.
- 874 52. Nawrocki, E. P. Structural RNA Homology Search and Alignment Using Covariance  
875 Models. PhD Thesis, Washington University in St. Louis, 2009.
- 876 53. Quast, C.; Pruesse, E.; Yilmaz, P.; Gerken, J.; Schweer, T.; Yarza, P.; Peplies, J.; Glöckner,  
877 F. O., The SILVA ribosomal RNA gene database project: improved data processing and  
878 web-based tools. *Nucleic Acids Research* **2012**, *41*, (D1), D590–D596.
- 879 54. Kanehisa, M.; Sato, Y.; Kawashima, M.; Furumichi, M.; Tanabe, M., KEGG as a reference  
880 resource for gene and protein annotation. *Nucleic Acids Research* **2016**, *44*, (D1), D457-  
881 D462.
- 882 55. Aramaki, T.; Blanc-Mathieu, R.; Endo, H.; Ohkubo, K.; Kanehisa, M.; Goto, S.; Ogata, H.,  
883 KofamKOALA: KEGG ortholog assignment based on profile HMM and adaptive score  
884 threshold. *bioRxiv* **2019**, 602110.
- 885 56. Eren, A. M.; Esen, Ö. C.; Quince, C.; Vineis, J. H.; Morrison, H. G.; Sogin, M. L.; Delmont,  
886 T. O., Anvi'o: an advanced analysis and visualization platform for 'omics data. *PeerJ* **2015**,  
887 *3*, e1319.
- 888 57. Alneberg, J.; Bjarnason, B. S.; de Bruijn, I.; Schirmer, M.; Quick, J.; Ijaz, U. Z.; Lahti, L.;  
889 Loman, N. J.; Andersson, A. F.; Quince, C., Binning metagenomic contigs by coverage and  
890 composition. *Nature Methods* **2014**, *11*, (11), 1144–1146.
- 891 58. Parks, D. H.; Imelfort, M.; Skennerton, C. T.; Hugenholtz, P.; Tyson, G. W., CheckM:  
892 assessing the quality of microbial genomes recovered from isolates, single cells, and  
893 metagenomes. *Genome Research* **2015**, *25*, (7), 1043–55.
- 894 59. Parks, D. H.; Rinke, C.; Chuvochina, M.; Chaumeil, P. A.; Woodcroft, B. J.; Evans, P. N.;  
895 Hugenholtz, P.; Tyson, G. W., Recovery of nearly 8,000 metagenome-assembled genomes  
896 substantially expands the tree of life. *Nature Microbiology* **2017**, *2*, (11), 1533–1542.
- 897 60. Olm, M. R.; Brown, C. T.; Brooks, B.; Banfield, J. F., DRep: A tool for fast and accurate  
898 genomic comparisons that enables improved genome recovery from metagenomes through  
899 de-replication. *ISME J.* **2017**, *11*, (12), 2864–2868.
- 900 61. Campbell, B. J.; Yu, L.; Heidelberg, J. F.; Kirchman, D. L., Activity of abundant and rare  
901 bacteria in a coastal ocean. *Proceedings of the National Academy of Sciences USA* **2011**,  
902 *108*, (31), 12776–12781.
- 903 62. Price, M. N.; Dehal, P. S.; Arkin, A. P., FastTree 2 - Approximately maximum-likelihood  
904 trees for large alignments. *PLoS One* **2010**, *5*, (3), e9490.

- 905 63. Rice, P.; Longden, L.; Bleasby, A., EMBOSS: The European Molecular Biology Open  
906 Software Suite. *Trends in Genetics* **2000**, *16*, (6), 276–277.
- 907 64. Quinlan, A. R.; Hall, I. M., BEDTools: a flexible suite of utilities for comparing genomic  
908 features. *Bioinformatics* **2010**, *26*, (6), 841-842.
- 909 65. Bushnell, B. BMap short-read aligner, and other bioinformatics tools.  
910 <http://sourceforge.net/projects/bbmap/>
- 911 66. Oksanen Blanchet, F. G., Kindt, R., Legendre, P., Minchin, P.R., O'Hara, R.B., Simpson,  
912 G.L., Solymos, P., Stevens, M.H.H., Wagner, H., J, *vegan: Community Ecology Package*.  
913 2013; Vol. R package
- 914 67. Love, M. I.; Huber, W.; Anders, S., Moderated estimation of fold change and dispersion for  
915 RNA-seq data with DESeq2. *Genome Biology* **2014**, *15*, (12), 550.
- 916 68. Volk, C.; Dundore, E.; Schiermann, J.; LeChevallier, M., Practical evaluation of iron  
917 corrosion control in a drinking water distribution system. *Water Research* **2000**, *34*, (6),  
918 1967-1974.
- 919 69. Santos, Q. M. B.-d. I.; L. Schroeder, J.; C. Sevillano-Rivera, M.; Sungthong, R.; Z. Ijaz, U.;  
920 T. Sloan, W.; J. Pinto, A., Emerging investigators series: microbial communities in full-  
921 scale drinking water distribution systems – a meta-analysis. *Environmental Science: Water*  
922 *Research & Technology* **2016**, *2*, (4), 631-644.
- 923 70. Berg, I. A., Ecological Aspects of the Distribution of Different Autotrophic CO<sub>2</sub> Fixation  
924 Pathways. *Applied and Environmental Microbiology* **2011**, *77*, (6), 1925-1936.
- 925 71. Caspi, R.; Billington, R.; Fulcher, C. A.; Keseler, I. M.; Kothari, A.; Krummenacker, M.;  
926 Latendresse, M.; Midford, P. E.; Ong, Q.; Ong, W. K.; Paley, S.; Subhraveti, P.; Karp, P. D.,  
927 The MetaCyc database of metabolic pathways and enzymes. *Nucleic Acids Research* **2018**,  
928 *46*, (D1), D633-D639.
- 929 72. Coleman, S. T.; Fang, T. K.; Rovinsky, S. A.; Turano, F. J.; Moye-Rowley, W. S.,  
930 Expression of a Glutamate Decarboxylase Homologue Is Required for Normal Oxidative  
931 Stress Tolerance in *Saccharomyces cerevisiae*. *Journal of Biological Chemistry* **2001**, *276*,  
932 (1), 244-250.
- 933 73. Feehily, C.; O'Byrne, C. P.; Karatzas, K. A. G., Functional  $\gamma$ -Aminobutyrate Shunt in  
934 *Listeria monocytogenes*: Role in Acid Tolerance and Succinate Biosynthesis. *Applied and*  
935 *Environmental Microbiology* **2013**, *79*, (1), 74-80.
- 936 74. Feehily, C.; Karatzas, K. a. G., Role of glutamate metabolism in bacterial responses towards  
937 acid and other stresses. *Journal of Applied Microbiology* **2013**, *114*, (1), 11-24.
- 938 75. Metzner, M.; Germer, J.; Hengge, R., Multiple stress signal integration in the regulation of  
939 the complex  $\sigma$ S-dependent *csiD-ygaF-gabDTP* operon in *Escherichia coli*. *Molecular*  
940 *Microbiology* **2004**, *51*, (3), 799-811.
- 941 76. Barberán, A.; Ramirez, K. S.; Leff, J. W.; Bradford, M. A.; Wall, D. H.; Fierer, N., Why are  
942 some microbes more ubiquitous than others? Predicting the habitat breadth of soil bacteria.  
943 *Ecology Letters* **2014**, *17*, (7), 794-802.
- 944 77. Guieysse, B.; Wuertz, S., Metabolically versatile large-genome prokaryotes. *Current*  
945 *Opinion in Biotechnology* **2012**, *23*, (3), 467-473.
- 946 78. Ahn, S.; Jung, J.; Jang, I.-A.; Madsen, E. L.; Park, W., Role of Glyoxylate Shunt in  
947 Oxidative Stress Response. *Journal of Biological Chemistry* **2016**, *291*, (22), 11928-11938.
- 948 79. Chatzigiannidou, I.; Props, R.; Boon, N., Drinking water bacterial communities exhibit  
949 specific and selective necrotrophic growth. *npj Clean Water* **2018**, *1*, (1), 1-4.

- 950 80. Dolan, S. K.; Welch, M., The Glyoxylate Shunt, 60 Years On. *Annual Review of*  
951 *Microbiology* **2018**, *72*, (1), 309-330.
- 952 81. Galhardo, R. S.; Rocha, R. P.; Marques, M. V.; Menck, C. F. M., An SOS-regulated operon  
953 involved in damage-inducible mutagenesis in *Caulobacter crescentus*. *Nucleic Acids*  
954 *Research* **2005**, *33*, (8), 2603-2614.
- 955 82. Baharoglu, Z.; Mazel, D., SOS, the formidable strategy of bacteria against aggressions.  
956 *FEMS Microbiology Reviews* **2014**, *38*, (6), 1126-1145.
- 957 83. Heal, K. R.; Qin, W.; Ribalet, F.; Bertagnolli, A. D.; Coyote-Maestas, W.; Hmelo, L. R.;  
958 Moffett, J. W.; Devol, A. H.; Armbrust, E. V.; Stahl, D. A.; Ingalls, A. E., Two distinct  
959 pools of B12 analogs reveal community interdependencies in the ocean. *Proceedings of the*  
960 *National Academy of Sciences USA* **2017**, *114*, (2), 364-369.
- 961 84. Chiao, T.-H.; Clancy, T. M.; Pinto, A.; Xi, C.; Raskin, L., Differential Resistance of  
962 Drinking Water Bacterial Populations to Monochloramine Disinfection. *Environmental*  
963 *Science & Technology* **2014**, *48*, (7), 4038-4047.