

1 **Evaluation of 16S next-generation sequencing of hypervariable**
2 **region 4 in wastewater samples: an unsuitable approach for**
3 **bacterial enteric pathogen identification**

4

5 Telleasha L Greay^{1,2}, Alexander W Gofton¹, Alireza Zahedi^{1,2}, Andrea Papparini¹, Kathryn L
6 Linge^{3,4}, Cynthia A Joll³ and Una M Ryan^{1*}

7

8 ¹Vector and Waterborne Pathogens Research Group, School of Veterinary and Life Sciences,
9 Murdoch University, Perth, Western Australia, Australia

10 ²Western Australian State Agricultural Biotechnology Centre, Murdoch University, Perth,
11 Western Australia, Australia

12 ³Curtin Water Quality Research Centre, Chemistry, School of Molecular and Life Sciences,
13 Curtin University, GPO Box U1987, Perth, Australia

14 ⁴ChemCentre, PO Box 1250, Perth, Western Australia, Australia

15

16 * Correspondence: una.ryan@murdoch.edu.au

17

18 Emails:

19 TLG: telleasha.greay@outlook.com

20 AWG: alexander.gofton@csiro.au

21 AZ: a.zahediabdi@murdoch.edu.au

22 AP: a.paparini@murdoch.edu.au

23 KLL: klinge@chemcentre.wa.gov.au

24 CAJ: c.joll@curtin.edu.au

25 **Abstract**

26 Recycled wastewater can carry human-infectious microbial pathogens and therefore
27 wastewater treatment strategies must effectively eliminate pathogens before recycled
28 wastewater is used to supplement drinking and agricultural water supplies. This study
29 characterised the bacterial composition of four wastewater treatment plants (WWTPs) (three
30 waste stabilisation ponds and one oxidation ditch WWTP using activated sludge treatment) in
31 Western Australia. The hypervariable region 4 (V4) of the bacterial 16S rRNA (16S) gene was
32 sequenced using next-generation sequencing (NGS) on the Illumina MiSeq platform.
33 Sequences were pre-processed in USEARCH v10.0 and denoised into zero-radius taxonomic
34 units (ZOTUs) with UNOISE3. Taxonomy was assigned to the ZOTUs using QIIME 2 and the
35 Greengenes database and cross-checked with the NCBI nr/nt database. Bacterial composition
36 of all WWTPs and treatment stages (influent, intermediate and effluent) were dominated by
37 Proteobacteria (29.0-87.4%), particularly Betaproteobacteria (9.0-53.5%) and
38 Gammaproteobacteria (8.6-34.6%). Nitrifying bacteria (*Nitrospira* spp.) were found only in
39 the intermediate and effluent of the oxidation ditch WWTP, and denitrifying and floc-forming
40 bacteria were detected in all WWTPs, particularly from the families Comamonadaceae and
41 Rhodocyclales. Twelve pathogens were assigned taxonomy by the Greengenes database, but
42 comparison of sequences from genera and families known to contain pathogens to the NCBI
43 nr/nt database showed that only three pathogens (*Arcobacter venerupis*, *Laribacter*
44 *hongkongensis* and *Neisseria canis*) could be identified in the dataset at the V4 region.
45 Importantly, Enterobacteriaceae genera could not be differentiated. Family level taxa assigned
46 by Greengenes database agreed with NCBI nr/nt in most cases, however, BLAST analyses
47 revealed erroneous taxa in Greengenes database. This study highlights the importance of
48 validating taxonomy of NGS sequences with databases such as NCBI nr/nt, and recommends
49 including the V3 region of 16S in future short amplicon NGS studies that aim to identify

50 bacterial enteric pathogens, as this will improve taxonomic resolution of most, but not all,
51 Enterobacteriaceae species.

52

53 **Keywords:** Wastewater, next-generation sequencing, 16S rRNA, V4, Greengenes,
54 Enterobacteriaceae.

55

56 **1. Introduction**

57 Water is becoming an increasingly scarce global resource, and as the overall demand
58 for water grows, the quantity of wastewater produced and its overall pollution load are
59 continuously increasing worldwide (Connor et al., 2017). Recycled wastewater is an essential
60 resource in addressing this problem, as properly treated water can be safely released back into
61 the environment, and used to supplement limited drinking water supplies. However, unless
62 effectively treated, recycled wastewater has the potential to carry microbial pathogens (viruses,
63 bacteria, protozoa and helminths), toxic chemicals and heavy metals. Therefore, treatment
64 strategies must effectively eliminate these major public health risks (Rodriguez-Manzano et
65 al., 2012).

66 Wastewater recycling in urban areas typically employs reverse osmosis membranes or
67 advanced oxidation treatment after activated sludge wastewater treatment. This results in high
68 purity recycled water, fit for potable reuse, but is technically challenging and expensive
69 (Rajasulochana and Preethy, 2016; Garrido-Cardenas et al., 2017). In contrast, rural WWTPs
70 typically use simple, non-mechanical waste stabilisation ponds (WSPs) or lagoons consisting
71 of open basins that rely on natural microorganisms and algae to assist in the breakdown and
72 settlement of degradable organic matter. Wastewater "influent" enters on one side of the WSP
73 and exits on the other side as "effluent", after spending days or even months undergoing
74 treatment processes in the pond, depending on plant capacity and flow rate. The treated effluent

75 is discharged generally for non-potable purposes, such as irrigation of public open spaces or
76 agricultural/horticultural uses (Von Sperling, 2007; Anon, 2009). These WSPs are widely used
77 across the world as passive wastewater treatment for domestic wastewaters as they can offer
78 low cost, low maintenance and effective pathogen removal (Von Sperling, 2007; Ho et al.,
79 2017; Eland et al., 2018).

80 Removal and inactivation of pathogens from WSPs is achieved via long retention times,
81 increased temperature and pH, the presence of algal antibacterial compounds and sunlight
82 penetration. Therefore shallow (<1 m) WSPs with low turbidity, high pH and maximal sunlight
83 exposure will achieve the most efficient pathogen removal (Sharafi et al., 2012). While WSP
84 systems can achieve high removal efficiencies (4-6 log₁₀), the efficiency of pathogen removal
85 in full-scale systems is highly variable, and many WSP systems achieve only 2-3 log₁₀ removal
86 (Verbyla et al., 2017).

87 In contrast to WSPs, many conventional WWTPs use an activated sludge process in
88 which a biological sludge containing living microorganisms is mixed with wastewater and
89 aerated in a reactor, forming a mixed liquor. Microbial populations within the activated sludge
90 include a range of bacteria, yeast, fungi, protozoa and higher organisms such as rotifers that
91 can digest organic matter in wastewater, and clump together (by flocculation), producing a
92 treated wastewater that is relatively free from suspended solids and organic material. The
93 removal mechanisms of pathogens in an activated sludge system are inactivation, hunting by
94 ciliate protozoa, adsorption to solids and capsulation inside the sludge flocs (Sharafi et al.,
95 2012).

96 Understanding the diversity of bacterial microorganisms in wastewater is essential for
97 understanding the performance for biological wastewater treatment systems (Inaba et al.,
98 2017). DNA-based approaches for identification of bacteria, such as polymerase chain reaction
99 (PCR) and Sanger sequencing, can overcome the limitations of conventional bacterial

100 identification techniques (e.g. microscopy, culture-dependent assays and biochemical
101 techniques) that are laborious and time-consuming, by allowing for the identification of
102 microbes that are morphologically indistinguishable, uncultivable, fastidious, and obligate
103 intracellular. Molecular bacterial identification approaches often target the 16S rRNA (16S)
104 gene, which enables species differentiation based on genetic dissimilarity. However, the
105 throughput of species identification with PCR and Sanger sequencing is limited by individual
106 clone library preparation, and species-specific PCR approaches require *a priori* hypotheses
107 regarding the taxa to be targeted. Wastewater can be comprised of hundreds of bacterial species
108 (Berlec 2012; Kim et al. 2015), therefore assessments of bacterial diversity on this scale using
109 PCR and/or Sanger sequencing is impractical. The rapid advances of next-generation
110 sequencing (NGS) technologies have revolutionised the ability to identify large numbers of
111 bacteria from various types of environmental and biological samples (Garrido-Cardenas et al.,
112 2017). Primers targeting one or more of the nine hypervariable (V) regions of 16S can be used
113 with NGS to identify bacteria. Other studies that have used 16S NGS to identify bacteria in
114 WWTPs have targeted V3-4 (Lu et al., 2015), V4 (Zhang et al., 2012) and V5-6 (McLellan et
115 al., 2010), and there is no consensus on the most suitable region to target for bacterial
116 assessments in WWTPs. The V4 region of 16S is commonly targeted in microbiome studies
117 with the widely used 515F/806R primers (Caporaso et al., 2011). These primers are
118 recommended in the Earth Microbiome Project's Illumina NGS protocol
119 (<http://www.earthmicrobiome.org/protocols-and-standards/16s/>) and have been modified by
120 other studies to include additional degeneracies to allow amplification of additional taxa
121 (Apprill et al., 2015; Parada et al., 2015). Therefore, the present study evaluated the ability of
122 the V4 16S NGS to identify bacteria, particularly enteric pathogens, in WWTPs, and used this
123 NGS approach to characterise bacterial compositions in different treatment stages (influent,
124 intermediate and effluent) of three WSPs and one oxidation ditch WWTP, which is a modified

125 activated sludge WWTP, that utilises prolonged aeration to remove biodegradable organic
126 compounds (Baars, 1962), in Western Australia (WA).

127

128 **2. Methods**

129 **2.1 Study sites and sample collection**

130 Wastewater samples ($n = 26$) were collected from three WSPs (WWTPs 1, 2 and 3) and
131 an oxidation ditch (WWTP 4) in 2015 in WA (Table 1 and Figure 1). Samples were collected
132 in February, July and September in 2015 and covered two seasons for each site. Samples were
133 collected from WWTP 1, located in north-west WA and in a tropical climate, during the wet
134 and dry seasons, while samples from WWTPs 2, 3 and 4, located in south-west WA and in a
135 temperate climate, were collected during summer and winter (Table 1). Wastewater samples
136 were also collected at different treatment stages (influent, intermediate and effluent) during
137 summer and winter (or dry and wet seasons for WWTP 1 samples) (Table 1). The wastewater
138 samples were collected in 1 L sterile containers that were treated with chlorine and rinsed with
139 the sample before filling. Samples were kept cool in an ice box during transport back to the
140 laboratory, and then stored at 4 °C and processed within 48 hours prior to DNA isolation.

141

142 **2.2 DNA isolation**

143 After 100 mL of each wastewater sample was filtered through sterile 0.2 µm Sterivex
144 filters (Millipore, USA), genomic DNA (gDNA) was extracted from the filters using a
145 PowerWater Sterivex DNA Isolation Kit (MO BIO Laboratories, California, USA). Extraction
146 reagent blank controls (ExCs; $n = 6$) were included alongside each batch of gDNA extractions.
147 Purified DNA was stored at -20 °C prior to molecular analysis.

148

149 **2.3 Next-generation sequencing library preparation**

150 The NGS library was prepared and sequenced following the 16S metagenomic
151 sequencing library preparation protocol from Illumina (Part # 15044223 Rev. B; Illumina,
152 USA), with minor modifications to the first stage PCR. V4 16S was amplified using modified
153 515F/806R primers [originally designed by Caporaso et al. (2011)]: 515FB 5'-
154 GTGYCAGCMGCCGCGGTAA-3' (Parada et al., 2015) and 806RB 5'-
155 GGACTACNVGGGTWTCTAAT-3' (Aprill et al., 2015). The 515FB/806RB primers were
156 modified to include Illumina MiSeq adapter sequences (Part # 15044223 Rev. B; Illumina,
157 USA), and conventional PCRs were carried out as described elsewhere
158 (www.earthmicrobiome.org/protocols-and-standards/16s/;
159 <https://doi.org/10.17504/protocols.io.nuudeww>). No-template controls (NTCs) were included
160 alongside each PCR. The V4 16S library was sequenced on the Illumina Miseq platform (San
161 Diego, CA, USA) with v2 sequencing chemistry.

162

163 **2.4 16S Bioinformatic analysis**

164 Paired-end 16S reads were merged (minimum 50 bp overlap), trimmed of primers and
165 distal bases, quality filtered (maximum expected error threshold of 1.0) and singletons were
166 removed with USEARCH v10.0 (Edgar, 2010), resulting in reads that were 247 bp in length
167 on average. Reads were denoised into zero-radius operational taxonomic units (ZOTUs) and
168 chimeras were filtered with UNOISE3 (Edgar, 2016). Taxonomic assignment of ZOTUs was
169 performed in QIIME 2 v2018.2 (Caporaso et al., 2010, <https://qiime2.org>) using the QIIME 2
170 feature classifier plugin (Bokulich et al, 2018) and the August 2013 release of the Greengenes
171 sequence database (McDonald et. al., 2012). The sequences were also BLAST searched against
172 the National Center for Biotechnology Information (NCBI) non-redundant nucleotide (nr/nt)
173 database to cross-check Greengenes assigned taxonomy. ZOTUs that were in low abundance

174 (<0.05% sequence composition) across all samples may represent PCR or sequencing error,
175 therefore, they were bioinformatically removed from the samples. To assess sequencing depth,
176 alpha rarefaction plots were generated with the R package vegan (Oksanen et al., 2018) using
177 R software (R Core Team, 2013).

178

179 **2.5 Phylogenetic analysis**

180 Enterobacteriaceae ZOTUs were aligned using the MAFFT program (Katoh et al.,
181 2002) with closely related sequences retrieved from the NCBI nr/nt database in Geneious
182 v10.2.2 (<http://www.geneious.com>, Kears e et al., 2012). Sequences in the alignment were
183 trimmed to the same length, then were imported into the program PhyML (Guindon et al.,
184 2010) and assessed for the most appropriate nucleotide substitution model (GTR+G+I) based
185 on Akaike Information Criterion (AIC). Maximum likelihood trees were constructed using
186 RAxML (Stamatakis, 2014). Genetic distance estimates were calculated with Kimura distance
187 matrices (Kimura, 1980) in Geneious v10.2.2.

188 ZOTU sequences generated from this study have been submitted to GenBank under the
189 accession numbers MH892609 to MH892828. Raw sequence files were deposited in the NCBI
190 Sequence Read Archive under the accession number PRJNA526519 (refer to Table 1 for
191 sample names and metadata).

192

193 **3. Results**

194 **3.1 Next-generation sequencing library summary**

195 Approximately 1.4 million paired-end V4 16S sequences were obtained for all samples
196 and controls ($n = 34$) (Table 2). After the reads were pre-processed (merged, quality filtered
197 with singletons and chimeras removed), there was a total of ~800,000 sequences for all samples
198 (~24,000 average). The processed 16S sequences (total of ~700,000) excluded sequences that

199 were not classified as bacteria and low abundance (<0.05%) ZOTUs, and on average, there
200 were ~27,000 processed bacterial 16S sequences for the WWTP samples ($n = 26$). Few
201 sequences were detected in the ExCs and NTCs, which had an average of 8 sequences (Table
202 2).

203

204 **3.2 Bacterial sequence composition in WWTPs**

205 A total of 3,598 ZOTUs (Supplementary File B.1) were obtained for the pre-processed
206 sequences, and a total of 1,644 ZOTUs remained for the processed sequences. For the
207 processed sequences, sequencing depth was adequate for all samples at ~5,000 sequences
208 (Supplementary Figures A.3 and A.4), but the alpha rarefaction plots did not reach a plateau
209 for the pre-processed sequences (Supplementary Figures A.1 and A.2). The archaeal sequence
210 compositions were low and two archaeal phyla were detected: Euryarchaeota was found in the
211 influent of WWTP 4 (<0.1%) and effluent of WWTP 2 (0.1%), and Parvarchaeota was found
212 in the effluent of WWTP 4 (0.1%). Two different types of Euryarchaeota were detected,
213 *Methanobrevibacter* sp. from the class Methanobacteria in WWTP 4 influent and
214 *Methanosaeta* sp. from the class Methanomicrobia in WWTP 2 effluent. The taxonomy for
215 Parvarchaeota was assigned as Parvarchaea for class, and WCHD3-30 and YLA114 for
216 Parvarchaea orders, with no further taxonomic classifications assigned by Greengenes
217 (Supplementary File B.2).

218 Bacteria were classified into 28 phyla (Supplementary File B.2); the most dominant
219 phylum was Proteobacteria across all WWTPs and treatment stages (influent, intermediate and
220 effluent), with sequence compositions ranging from 29.0% in the effluent of WWTP 2 to 87.4%
221 in the intermediate stage of WWTP 3 (Figure 2). Other abundant phyla (>10% composition in
222 WWTP samples) were Bacteroidetes (ranging from 4.1% in WWTP 1 influent to 31.5%
223 WWTP 3 effluent), Cyanobacteria (0% (not detected) in WWTP 1 and 3 influent to 47.2%

224 WWTP 2 effluent), Firmicutes (0.1% in WWTP 3 effluent to 22.1% in WWTP 1 influent) and
225 Actinobacteria (1.1% in WWTP 4 influent to 10.3% in WWTP 2 influent) (Figure 2).

226 Six classes of Proteobacteria were identified: Alphaproteobacteria, Betaproteobacteria,
227 Deltaproteobacteria, Epsilonproteobacteria, Gammaproteobacteria and "TA18".
228 Betaproteobacteria and Gammaproteobacteria sequences were abundant ($\geq 8.6\%$) and prevalent
229 across all WWTPs and treatment stages. There were also six classes for Bacteroidetes: WWTP
230 1 and 4 exhibited a similar pattern in sequence composition of Bacteroidetes, with classes
231 Bacteroidia and Flavobacteriia detected in the influent, and in addition to these two classes,
232 three other classes (Saprospirae and Sphingobacteriia) were also detected in the intermediate
233 and effluent stages (Figure 2). Like WWTP 1 and 4, Bacteroidia and Flavobacteriia were
234 detected in all stages of WWTP 3, and the same classes that were found in WWTP 1 and 4
235 were also found in WWTP 3, but in the intermediate stage of WWTP 3, only Bacteroidia,
236 Flavobacteriia and Saprospirae sequences were obtained. WWTP 2 had similar Bacteroidetes
237 in the influent and effluent; all aforementioned Bacteroidetes classes were found in the influent
238 and effluent, and an additional class, Rhodothermi, was also found in the effluent (Figure 2 and
239 Supplementary File B.2).

240 Cyanobacteria were not found in the influent of WWTP 1 and 3, but sequences were
241 detected in the intermediate and effluent stages of these plants, and were detected in all stages
242 of WWTP 2 and 4. Oscillatoriothycidae was dominant in the intermediate and effluent of
243 WWTP 1 and 2 (11.2% and 14.3%, respectively) and was also detected in the effluent of
244 WWTP 2 and 3. Other classes of Cyanobacteria included Synechococcophycidae in the
245 intermediate of WWTP 1 and effluent of WWTP 1 and 4, Nostocophycidae in WWTP 2
246 effluent, and a class designated as 4C0d-2 by the Greengenes database was found in WWTP 3
247 intermediate and WWTP 4 intermediate and effluent. Among the Firmicutes, three classes were
248 detected: Bacilli (Bacillales, Lactobacillales and Turicibacterales), Clostridia (Clostridiales)

249 and Erysipelotrichi (Erysipelotrichales). Bacilli and Clostridia sequences were the most
250 abundant classes of Firmicutes and were found in the influent of WWTPs 1, 3 and 4, ranging
251 from 0.8-22.1%, and sequences were in low abundance ($\leq 2.0\%$) or not detected in the
252 intermediate and effluent stages (Figure 2 and Supplementary File B.2).

253 Five Actinobacteria classes were identified: Acidimicrobiia, Actinobacteria,
254 Coriobacteriia, Nitriliruptoria and Thermoleophilia. The sequence composition of the class
255 Actinobacteria was higher than other classes of the phylum Actinobacteria (which were all
256 $\leq 2.0\%$ in various treatment plants and stages), particularly in the intermediate and effluent of
257 WWTPs 1, 3 and 4 (1.7-8.0%), and the influent and effluent of WWTP 2 (9.1% and 4.8%,
258 respectively) (Figure 2). Acidimicrobiia and Thermoleophilia were detected in the intermediate
259 and effluent of WWTP 1, influent and effluent of WWTP 2 and intermediate of WWTP 4. A
260 low sequence composition of Actinobacteria and Coriobacteriia were found in the influent of
261 WWTPs 1, 3 and 4. Nitriliruptoria was only found in the effluent of WWTP 2 (Supplementary
262 File B.2).

263

264 **3.3 Bacterial pathogen identification**

265 Based on Greengenes taxonomic assignments, seven ZOTUs were assigned to the
266 family Enterobacteriaceae (Gammaproteobacteria: Enterobacteriales). Four of these ZOTUs
267 were not assigned further taxonomy by Greengenes, but the remaining ZOTUs were designated
268 as *Citrobacter* sp., *Escherichia coli* and *Trabulsiella* sp. with high confidence (0.95-1).
269 However, comparison of the Enterobacteriaceae sp. ZOTUs to the NCBI nr/nt database using
270 BLAST revealed that all these ZOTUs were 100% similar to multiple Enterobacteriaceae sp.
271 genera (Table 3). The phylogenetic tree constructed with Enterobacteriaceae sp. ZOTUs and
272 Enterobacteriaceae sequences from the NCBI nr/nt database showed that different genera
273 grouped closely with short branch lengths, and most bootstrap values were low (Figure 3; refer

274 to Supplementary file B.3 for pairwise genetic distances, which range from 94.3-100%). This
275 supports the BLAST results from Table 3 that suggest that the Enterobacteriaceae sp. ZOTUs
276 can only be confidently assigned to the family level, and suggests that the V4 region of 16S
277 cannot distinguish between many Enterobacteriaceae species and genera.

278 Other ZOTUs from the class Gammaproteobacteria that were assigned to pathogenic
279 species, or to taxa that contain pathogens, based on Greengenes taxonomy included
280 *Acinetobacter* (*Acinetobacter johnsonii*, *Acinetobacter lwoffii* and unassigned species),
281 Aeromonadaceae (unassigned genera and *Tolumonas*), Coxiellaceae (genus unassigned),
282 Legionellaceae (genus unassigned), *Enterococcus* spp. (Enterococcaceae), Pseudomonadaceae
283 (*Pseudomonas alcaligenes*, *Pseudomonas fragi*, *Pseudomonas nitroreducens*, *Pseudomonas*
284 *stutzeri*, *Pseudomonas veronii*, *Pseudomonas viridiflava* and unassigned species),
285 Piscirickettsiaceae (genus unassigned) and Pseudoalteromonadaceae (genus unassigned).
286 Greengenes taxonomy that conflicted with BLAST analysis was identified for *Tolumonas*
287 (Aeromonadaceae; ZOTU 483), which was most similar to *Pseudaeromonas sharmana* (100%;
288 GenBank® accession no. MF280154), *Aeromonas sharmana* (99.2%; JF496528) and
289 *Tolumonas* sp. (98.8%; MG801837). *Acinetobacter* ZOTUs assigned to the species level
290 (*Acinetobacter johnsonii* and *Acinetobacter lwoffii*) with high confidence naïve Bayes
291 confidence scores (0.94-1) were also 100% similar to several other *Acinetobacter* species.
292 Similarly, many *Pseudomonas* species had sequence similarities of 100%, therefore had
293 incorrect species level taxonomy assigned with Greengenes. Greengenes taxonomy was more
294 conservative for ZOTU 589 (Pseudoalteromonadaceae sp.) as BLAST results showed that this
295 sequence could be assigned to the genus *Vibrio*, but like *Acinetobacter* and *Pseudomonas*,
296 many *Vibrio* species were also 100% similar at the V4 region (Supplementary File B.4).

297 The BLAST results agreed with Greengenes taxonomic assignments for most
298 Betaproteobacteria (Alcaligenaceae spp., Neisseriaceae spp. and *Vitreoscilla* spp.), except for

299 ZOTU 55 that was classed as *Microvirgula* sp., but was 100% similar to *Laribacter*
300 *hongkongensis* sequences (NR025167). Five Campylobacteraceae (class
301 Epsilonproteobacteria) ZOTUs that were assigned to the genus *Arcobacter* or *Arcobacter*
302 *cryaerophilus* were also 100% similar to *Campylobacter* sequences and therefore could only
303 be confidently assigned to the Campylobacteraceae family. *Corynebacterium* spp. and
304 *Mycobacterium* spp. (phylum Actinobacteria) taxonomy agreed for both Greengenes and
305 BLAST results, but there were discrepancies for Streptococcaceae spp. (phylum Firmicutes)
306 that were assigned as *Streptococcus* spp., *Streptococcus luteciae* and *Streptococcus minor* by
307 Greengenes, but could not be assigned to the species or genus level by BLAST in most cases.
308 The Greengenes taxa *Candidatus Rhabdochlamydia* sp. (ZOTU 1597; phylum Chlamydiae),
309 *Clostridium* spp., *Proteiniclasticum* sp. (phylum Firmicutes) or *Treponema* spp. (Spirochaetes)
310 could also not be confidently assigned to the genus level based on BLAST results
311 (Supplementary File B.4).

312 The sequence compositions for pathogenic and potentially pathogenic taxa that were
313 given final taxonomic assignments based on Greengenes and NCBI nr/nt sequence database
314 comparisons are summarised in Table 4. Briefly, sequence compositions for these taxa were
315 generally higher in the influent for WWTP 1, 3 and 4 and lower in the intermediate and effluent,
316 with the exception of *Acinetobacter* spp. in the intermediate stage of WWTP 3 (17.1%) and
317 *Aeromonas* sp. in the effluent of WWTP 1 (8.6%). Potentially pathogenic sequence
318 compositions were relatively low in the influent and effluent of WWTP 2, with the highest
319 composition of 2.0% in the effluent for Alcaligenaceae sp. (Table 4).

320

321 **3.4 Nitrifying, denitrifying and floc-forming bacteria**

322 Other bacteria of interest in WWTPs, such as nitrifying, denitrifying and floc-forming bacteria,
323 also had Greengenes taxonomy validated with BLAST results from the NCBI nr/nt database.

324 Compared to pathogenic bacteria, the nitrifying, denitrifying and floc-forming bacterial
325 ZOTUs had more taxonomic assignments that agreed with both databases. All were assigned
326 to the appropriate family, but some ZOTUs had conflicting genera. For example, ZOTU 387
327 was assigned as *Dechloromonas* sp. by Greengenes, but was also 100% similar to *Azonexus*
328 *hydrophilus* (LN650477), and ZOTU 766 Greengenes taxonomy was *Comamonas* sp., but this
329 ZOTU was 100% similar to *Comamonas* spp. (MH174324) and *Delftia* spp. (MF156914).
330 Results of taxonomy database comparisons for nitrifying, denitrifying and floc-forming
331 bacteria are provided in Supplementary File B.5, and the sequence compositions of validated
332 taxa are presented in Table 5. Nitrifying bacteria, *Nitrospira* spp. (Nitrospirales:
333 Nitrospiraceae), were only detected in the intermediate and effluent of WWTP 4, with sequence
334 compositions of 1.2% and 1.5%, respectively. In WWTP 1, denitrifying bacteria with the
335 highest compositions were found in the influent for *Comamonas* sp. (Comamonadaceae; 6.9%)
336 and *Thauera* spp. (Rhodocyclaceae; 3.4%). The comamonads *Hydrogenophaga* spp. and
337 *Aquabacterium* sp. had highest compositions in the effluent (2.4%) and influent (1.3%) of
338 WWTP 2, respectively, and the floc-forming bacteria *Flavobacterium* spp. were higher in the
339 effluent (4.7%) than in the influent (2.8%) of WWTP 2. WWTP 3 had a greater diversity of
340 denitrifying and floc-forming bacteria in the influent and intermediate stages than the effluent;
341 the highest sequence compositions in the influent was 4.4% for *Comamonas* sp., 6.0% for
342 *Thauera* spp. in the intermediate stage and 7.6% for *Flavobacterium* spp. in the effluent. The
343 abundance of *Comamonas* sp. sequences was also high in the influent of WWTP 4, and the
344 denitrifying bacteria *Uliginosibacterium* spp. were highest in the intermediate stage, and
345 *Flavobacterium* spp. was highest in the effluent of WWTP 4 (Table 5). Pseudomonadaceae
346 (*Pseudomonas*) are also denitrifying bacteria, and are summarised in Table 4.
347

348 **4. Discussion**

349 Evaluation of BLAST results from the NCBI nr/nt database of V4 16S sequences that were
350 assigned taxonomy by the 16S Greengenes taxonomy database to pathogenic species (or to
351 bacterial groups that contain pathogenic species) showed that the V4 region of 16S resolves
352 poorly at the species level, and genus level identification was also impeded in many instances.
353 Comparison of the ZOTU sequences to the NCBI nr/nt database revealed that only three
354 ZOTUs were 100% identical to the following pathogenic species: *Laribacter hongkongensis*,
355 which causes gastroenteritis and diarrhoea (Beilfuss et al., 2015); *Neisseria canis*, which
356 usually infects cats and dogs, but can also infect humans (Safton et al., 1999); and *Arcobacter*
357 *venerupis*. There are 15 species of *Arcobacter*, and three (*Arcobacter butzleri*, *Arcobacter*
358 *cryaerophilus* and *Arcobacter skirrowii*) have been associated with gastrointestinal infections
359 (Kayman et al., 2012). *Arcobacter venerupis* has previously only been isolated from shellfish
360 (Levican et al., 2012), and these sequences were in low abundance (0.7%) and only found in
361 the influent of WWTP 1. The sequences from *L. hongkongensis* and *N. canis* were found in the
362 influent of WWTPs 1, 3 and 4, and were in low abundance ($\leq 0.1\%$) or not detected in the
363 intermediate and effluent stages of these plants. Genera known to contain pathogenic species
364 that were validated by BLAST analyses of the ZOTUs against the NCBI nr/nt database
365 included *Aeromonas* sp., *Acinetobacter* spp., *Arcobacter* spp., *Candidatus Rhabdochlamydia*
366 sp., *Corynebacterium* sp., *Enterococcus* spp., *Legionella* sp., *Mycobacterium* spp., *Neisseria*
367 sp., *Pseudomonas* spp., *Streptococcus* spp., *Turneriella* sp. and *Vibrio* sp. A previous 16S NGS
368 study on WWTPs in Australia that also used the Illumina MiSeq platform identified 25
369 potentially pathogenic genera (Ahmed et al., 2017), while another study of municipal activated
370 sludge plants across four countries (China, USA, Canada and Singapore) identified 16
371 pathogenic genera using pyrosequencing (Ye and Zhang, 2011). The abundance of pathogenic
372 genera may vary among studies due to DNA extraction kits, different sequencing technologies,

373 inherent amplification biases during PCR and the 16S hypervariable region(s) targeted (Haft
374 and Tovchigrechko 2012).

375 Other pathogenic genera that can infect people via contaminated drinking water,
376 *Campylobacter* spp. and *Leptospira* spp., were not identified in the present study, but we cannot
377 exclude the possibility of their presence, as several *Campylobacteraceae* spp. and
378 *Leptospiraceae* spp. ZOTUs could not be resolved to the genus level (Table 4 and
379 Supplementary File B.4). Most of the potentially pathogenic genera identified had higher
380 sequence compositions in the influent, and had low compositions or were not detected in the
381 effluent (Table 4). However, *Aeromonas* sp. had relatively high sequence compositions in
382 WWTP 1 intermediate (2.2%) and effluent (8.6%) samples, and a similar trend was observed
383 for *Aeromonas* sp. in WWTP 3, with compositions of 4.7% in the influent, 1.1% in intermediate
384 samples and 4.8% in the effluent. *Acinetobacter* spp. also had a high sequence composition in
385 the intermediate samples of WWTP 3 (17.1%), but was not detected in the effluent. Other
386 studies have also found that *Acinetobacter* spp. sequence compositions were not significantly
387 lower in treated wastewater samples compared to the influent (Ahmed et al., 2017).
388 *Mycobacterium* spp. and *Pseudomonas* spp., which had lower compositions or were not
389 detected in the influent, had higher compositions in the intermediate and effluent (Table 4).
390 The absence or lower abundance of bacteria associated with human waste in the influent
391 compared to the intermediate and effluent may be partly explained by the lack of sample
392 replicates, as only 100 mL grab samples were collected per season at each site and treatment
393 stage. However, the number of 16S sequences obtained by NGS does not represent the number
394 of bacterial organisms present. A number of factors affect sequence composition, including
395 PCR amplification bias (Hong et al., 2009), sequencing depth and copy number variation in
396 the 16S gene (Kembel et al., 2012).

397 Many enteric bacteria (Enterobacteriaceae) can be transmitted to humans by faecal-oral
398 transmission and can cause gastrointestinal illnesses with symptoms of abdominal pain,
399 diarrhoea, fever, nausea and vomiting. Human enteric pathogens include *Citrobacter freundii*,
400 *Escherichia coli*, *Klebsiella aerogene*, *Salmonella bongori*, *Salmonella enterica* and *Shigella*
401 spp. (Cabral, 2010). Other pathogens from the family Enterobacteriaceae that can cause
402 gastrointestinal illnesses are *Yersinia enterocolitica*, which is a food-borne pathogen associated
403 with pork products (Bhaduri et al., 2005), and *Raoultella ornithinolytica* (formerly *Klebsiella*
404 *ornithinolytica*), which has been found in aquatic environments and hospitals, with one report
405 of its isolation from human digestive organs (Seng et al., 2016). Enterobacteriaceae pathogens
406 that cause urinary tract infections and other illnesses in humans include *Proteus mirabilis*,
407 *Proteus penneri*, *Proteus vulgaris* and *Serratia marcescens* (Guentzel et al., 1996).
408 Unfortunately, the V4 region of 16S lacked sufficient variability to distinguish between
409 Enterobacteriaceae genera (Table 3 and Figure 3). The same issue was likely encountered in
410 the V4 16S NGS study by Zhang et al. (2012), that reported the detection of sequences from
411 the order Enterobacteriales. Similarly, a 16S NGS study on WWTPs that targeted the V6 region
412 could not resolve one OTU that was 100% similar to several Enterobacteriaceae genera
413 (primarily *Klebsiella* and *Shigella*) (McLellan et al., 2010). Other 16S wastewater studies that
414 have targeted 16S regions that span two hypervariable regions appear to have been able to
415 resolve Enterobacteriaceae genera. For example, Ahmed et al. (2017) sequenced regions V5-6
416 (300 bp), and reported the detection of *Escherichia/Shigella* (unclear if these could be
417 differentiated), *Salmonella* and *Yersinia*, but species level assignments were not made. Lu et
418 al. (2015) targeted the V3-4 region (460 bp) and reported the presence of *Klebsiella*
419 *pneumoniae* and *Serratia* spp., but performed shotgun sequencing to identify pathogens to the
420 species level, which included *E. coli*, *S. enterica*, *Shigella sonnei* and *Yersinia pestis*.
421 According to a study by Chakravorty et al. (2007), V3 is a more suitable region for the

422 differentiation of Enterobacteriaceae genera, and these authors recommended targeting V2, V3
423 and V6 to identify the bacterial genera assessed in their study, including *Acinetobacter*,
424 *Bacillus*, *Clostridium*, *Corynebacterium*, *Chlamydia*, *Enterococcus*, *Listeria*, *Mycobacterium*,
425 *Neisseria*, *Pseudomonas*, *Streptococcus*, *Staphylococcus*, *Treponema* and *Vibrio*. Using these
426 three regions means that most of the 110 species examined in their study could be identified to
427 the species level (Chakravorty et al., 2007). Using multiple regions does have some challenges,
428 however. For example, the V2 region of *E. coli* starts at nucleotide (nt) position 137 and V6
429 ends at nt position 1,043 (Brosius et al., 1978), therefore V2-6 spans 906 bp of 16S. This
430 amplicon is too long for current amplicon NGS sequencers; the maximum length is 600 bp on
431 the Illumina MiSeq with v3 chemistry (<http://www.illumina.com/>). Regions V2-3 and V6
432 could be targeted separately, or full length 16S could be sequenced on long-read sequencing
433 platforms such as PacBio for improved taxonomic resolution of a greater variety of taxa (Ibal
434 et al. 2019). It is important for Enterobacteriaceae species such as *Escherichia coli* for
435 serotypes to be differentiated at the strain level, as some strains are harmless gut bacteria
436 whereas others are pathogenic, e.g. enterohemorrhagic *Escherichia coli* O157:H7. While some
437 studies state that 16S sequencing is unsuitable for differentiating *E. coli* and *Shigella* spp.
438 serotypes as the sequence similarity is high (97.9-99.9%) (Devanga Ragupathi et al. 2017), Ibal
439 et al. (2019) were able to classify *E. coli* strains based on full length 16S sequences (Ibal et al.
440 2019). Other housekeeping genes that are conserved among bacteria, such as *gyrB*, *rpoB*
441 and *mdh* have greater genetic variability for distinguishing *E. coli* and *Shigella* spp. strains than
442 16S (Devanga Ragupathi et al. 2017; Fukushima et al. 2002). These genes could also be
443 targeted using amplicon NGS approaches for improved taxonomic resolution of bacterial
444 strains, however the use of universal primers is more limited than 16S. Alternatively, shotgun
445 sequencing could be performed, which can provide greater taxonomic and functional
446 information (e.g. pathogenicity islands and toxin-producing genes) than amplicon NGS of

447 several target genes (Sanapareddy et al., 2009; Lu et al., 2015). Shotgun sequencing of
448 metagenomes has been considerably more expensive than amplicon NGS (Goodwin et al.,
449 2016), however costs are reducing, particularly with new approaches such as "shallow shotgun
450 sequencing", which can produce more accurate species level taxonomic and functional profiles
451 of the human microbiome than 16S sequencing (Hillmann et al. 2018).

452 A large portion of the V4 16S sequences (68%) collected in this current study were not
453 assigned to the genus level with the Greengenes database. Other 16S NGS studies on
454 wastewater have used RDP Classifier (Zhang et al., 2012; Ahmed et al., 2017) and SILVA
455 (McLellan et al., 2010; Lu et al., 2015) databases for taxonomic assignment. According to a
456 recent study that compared the major taxonomy databases (Greengenes, RDP classifier,
457 SILVA, NCBI and OTT), there were few conflicts when SILVA, RDP and Greengenes were
458 mapped into NCBI and OTT (Balvočiūtė et al., 2017). However, we found many genus level
459 conflicts, when potentially pathogenic and denitrifying bacteria were compared to the NCBI
460 nr/nt database (Supplementary Files B.4 and B.5). Furthermore, we found erroneous taxonomy
461 in the Greengenes database that causes 16S sequences deriving from chloroplasts in algae and
462 plants to be classified to the bacterial phylum Cyanobacteria and the class "Chloroplast", which
463 is not a valid taxon. For 44 ZOTUs in our dataset that were classified to the class "Chloroplast",
464 the orders provided by the Greengenes database were Chlorophyta (phylum of green algae),
465 Euglenozoa (phylum of flagellate excavates) and Stramenopiles (infrakingdom of algae and
466 oomycetes). While chloroplast sequences in the Greengenes database can be useful to identify
467 such sequences in an NGS dataset, researchers that aim to only analyse bacterial 16S sequences
468 at higher levels of classification (kingdom and phylum) may be unaware that the chloroplast
469 sequences are classified in the database at the class level. Classifying the chloroplast sequences
470 as "Chloroplast" at the kingdom level, rather than as "Bacteria" may help researchers to detect
471 these sequences at an earlier stage of the data analysis. We have provided a modified version

472 of the Greengenes 99 OTU taxonomy file for all chloroplast sequences with the kingdom
473 “Bacteria” renamed as “Chloroplast” in Supplementary File B.6. A custom curated sequence
474 database for waterborne pathogens, with quality-checked sequences and taxonomy validated
475 by phylogenetic analyses, may also reduce the errors in bacterial taxonomic assignment
476 experienced with other 16S sequence databases.

477 Overall, of the 2 archaeal and 28 bacterial phyla detected, Proteobacteria, Bacteroidetes,
478 Cyanobacteria, Firmicutes and Actinobacteria had high sequence compositions (>10%) in
479 WWTP samples (Figure 2). The two most dominant phyla in all treatment stages for WWTPs
480 1-4 were Proteobacteria and Bacteroidetes, which has also been observed by a previous 16S
481 NGS study that examined bacteria in activated sludge WWTPs across Australia, including
482 Perth (Ahmed et al., 2017). The study by McLellan et al. (2010) that compared V6 16S NGS
483 bacterial profiles in WWTP influent, surface water and human faecal samples, also found that
484 the most dominant bacterial phylum in the WWTPs was Proteobacteria (overall 59% sequence
485 composition), and like our study, Gammaproteobacteria and Betaproteobacteria were the most
486 abundant classes. McLellan et al. (2010) also found that Actinobacteria, Bacteroidetes and
487 Firmicutes were dominant taxa in the WWTP influent, and sewage samples had high
488 compositions of Firmicutes, particularly Clostridia (the human faecal samples were comprised
489 mostly (98%) of Clostridia) and Bacilli (McLellan et al., 2010). In the present study, Firmicutes
490 had the highest compositions in the influent of WWTPs 1, 3 and 4, ranging from 16.8-20.5%;
491 Bacilli ranged from 5.6-11.9% and Clostridia ranged from 7.9-12.4% (Figure 2). *Bacteroides*
492 is another faecal indicator bacterium (Kreader, 1995), and the sequence compositions for
493 *Bacteroides* spp. in the present study ranged from 0.4% in the influent of WWTP 1 and 2 to
494 2.4% in the influent of WWTP 3, and sequence compositions were low ($\leq 0.8\%$) or undetectable
495 in the intermediate and effluent (Supplementary file B.2). *Faecalibacterium* is also associated
496 with faeces (Zheng et al., 2009), and was detected in the influent of WWTP 1 (1.0%), 3 (1.7%)

497 and 4 (1.3%), and had low sequence compositions ($\leq 0.1\%$) or were not detected in the
498 intermediate and effluent stages.

499 Nitrification is a fundamental process in the biological removal of nitrogen in WWTPs,
500 and this two-step process is carried out by ammonia-oxidising bacteria (AOB) that convert
501 ammonia to nitrite, then nitrite-oxidising bacteria (NOB) convert nitrite to nitrate (Bellucci and
502 Curtis, 2011). *Nitrosomonas* and *Nitrospira* are two important genera of AOB in WWTPs,
503 while *Nitrobacter* is a major NOB (Siripong et al., 2007). In the present study, *Nitrosomonas*
504 and *Nitrobacter* were not detected, and *Nitrospira* spp. were only detected in the intermediate
505 and effluent of WWTP 4 (sequence compositions 1.2% and 1.5%, respectively) (Table 5).
506 Rhodocyclales are a widespread and abundant order of bacteria in WWTPs responsible for
507 anaerobic nitrogen removal by denitrification (Yang et al., 2011). In the present study, 12
508 Rhodocyclales genera were identified: *Azoarcus* spp., *Azonexus* spp., *Azospira* spp.,
509 *Dechloromonas* spp., *Methyloversatilis* sp., *Propionivibrio* spp., *Rhodocyclaceae* spp.,
510 *Sterolibacterium* spp., *Sulfuritalea* sp., *Thauera* spp., *Uliginosibacterium* spp. and *Zoogloea*
511 spp. (Table 5). In WWTP 1, Rhodocyclales were highest in the influent (*Thauera* spp. had the
512 highest composition; 3.4%) and rare ($\leq 0.1\%$) or not detected in the intermediate and effluent
513 samples. Rhodocyclales compositions were low in the influent (0.8%) and effluent (1.0%) of
514 WWTP 2. WWTP 3 had higher Rhodocyclales compositions in the intermediate (13.2%; most
515 abundant was *Thauera* spp. at 6.0%) compared to the influent (4.3%), and no Rhodocyclales
516 sequences were detected in WWTP 3 effluent. In addition to denitrification, certain *Thauera*
517 and *Dechloromonas* strains can degrade oil derivatives such as toluene (Shinoda et al., 2004;
518 Chakraborty et al., 2005) and therefore may be important in reducing the ecological burden of
519 these aromatic compounds, but we were unable to identify the species and strains of these
520 genera based on V4 16S amplicons. Unlike the WSPs, the oxidation ditch plant WWTP 4 had
521 high Rhodocyclales compositions in both the intermediate and effluent (7.7% and 7.1%,

522 respectively). Members of the family Comamonadaceae are also denitrifiers and are
523 responsible for aromatic degrading processes (Xu et al., 2018). The nine Comamonadaceae
524 genera identified were *Aquabacterium* sp., *Brachymonas* (*Brachymonas denitrificans*),
525 *Comamonas* sp., *Delftia* sp., *Flavobacterium* spp., *Hydrogenophaga* spp., *Polaromonas* spp.,
526 *Rhodoferax* spp. and *Rubrivivax* spp. Comamonadaceae compositions in WWTP 1 were similar
527 to those observed for Rhodocyclales in this treatment plant, with the highest composition
528 observed in the influent (7.4%; *Comamonas* sp. had the highest composition of 6.9%) and
529 compositions were low in the intermediate and effluent (1.6% and 1.7%, respectively).
530 Comamonadaceae compositions were much higher than Rhodocyclales in WWTP 2, which
531 had 4.7% in the influent and 7.7% in the effluent, and *Flavobacterium* spp. had the greatest
532 sequence compositions in both influent (2.8%) and effluent (4.7%). For WWTP 3, the
533 Comamonadaceae compositions were similar to Rhodocyclales in the influent and
534 intermediate, but unlike Rhodocyclales, were detected (mostly *Flavobacterium* spp. 7.6%) in
535 the effluent. For WWTP 4, the Comamonadaceae were mostly comprised of *Comamonas* sp.
536 in the influent (4.4%) and *Flavobacterium* spp. (6.5%) in the effluent, and the composition of
537 Comamonadaceae was low in the intermediate (1.2%). Comamonadaceae, Rhodocyclaceae,
538 Flavobacteriaceae and Pseudomonadaceae also play important roles in flocculation in activated
539 sludge plants, and Comamonadaceae and Flavobacteriaceae are important for bulking and
540 foaming (Shchegolkova et al., 2016).

541 **5. Conclusions**

542 In the present study, a total of 36 pathogenic or potentially pathogenic species were detected,
543 but most could not be identified to species level. Of these, sequences belonging to 14 medically
544 important genera that could possibly be from pathogens were identified primarily in the
545 influent of WWTPs 1-4. In almost all cases, these bacteria were present in lower abundance in
546 the effluent with the exception of *Aeromonas* sp. in the effluent of WWTP 1 (8.6%). The use

547 of V4 16S NGS for bacterial pathogen identification has significant limitations for species level
548 identification including the inability to differentiate Enterobacteriaceae genera that contain
549 many important enteric pathogens of humans. Amplicon NGS is a useful tool for broad
550 taxonomic surveys of bacteria, while tools such as quantitative PCR and droplet digital PCR
551 could be used in follow-up studies to identify bacteria that could not be differentiated at the
552 species or strain level. This would also allow quantification of pathogens before and after the
553 wastewater treatment process. Future studies that aim for improved taxonomic resolution of
554 bacterial pathogens in wastewater should consider sequencing full length 16S and more
555 variable housekeeping genes such as *gyrB*, *rpoB* or *mdh* for differentiation of *E. coli* and
556 *Shigella* strains. Shallow shotgun sequencing can also be used for pathogen identification and
557 for gaining functional information that is important for public health.

558 Nitrifying, denitrifying and floc-forming bacteria could mostly be identified to the
559 genus level. Only the activated sludge oxidation ditch plant showed the presence of an AOB,
560 *Nitrospira* spp., for bacterial nitrification. However, both the lower technology WSPs and the
561 activated sludge oxidation ditch plant showed the presence of Rhodocyclales,
562 Comamonadaceae, Flavobacteriaceae and Pseudomonadaceae bacteria, which are responsible
563 for anaerobic nitrogen removal by denitrification (i.e. conversion of nitrate to nitrogen gas).
564 These bacteria are also important for WWTP performance since they assist floc formation. Our
565 current work is examining the presence, diversity and relative abundances of bacterial
566 communities responsible for the nitrification and denitrification cycle in WSPs
567 (e.g. *Nitrobacter*, *Nitrosomonas*, *Nitrospira*, *Nitrosococcus* and *Nitrosomonas*) using
568 functional genes that encode key enzymes (*amoA*, *njfH*, *nirK*, *nosZ*, *norB*, *nxB*, *narG*, *napA*
569 and *nrfA*). This will help us to better understand the correlations between the concentrations of
570 selected nitrogenous species present in wastewater and their contribution to the nitrogen cycle
571 in WSPs.

572 Other limitations include the misidentification of 16S sequences from chloroplasts as
573 Cyanobacteria by the Greengenes database. Due to the discrepancies between taxonomic
574 assignments with Greengenes and the NCBI nr/nt database, we recommend that future studies
575 use the Greengenes database for 16S NGS taxonomic assignment with caution and compare
576 OTU or ZOTU sequences with the NCBI nr/nt database to validate taxonomic assignments.

577

578 **Acknowledgements**

579 This work was supported by the Australian Research Council (ARC LP130100602), in
580 collaboration with the Water Corporation of Western Australia and Water Research Australia.
581 We thank Dr Elvina Lee, Annachiara Codello and Maninder Khurana for their assistance with
582 NGS preparation. We thank Arron Lethorn and operational staff at the Water Corporation of
583 Western Australia for project management and assisting with site visits and sample collection.

584

585 **References**

- 586 Ahmed W., Staley C., Sidhu J., Sadowsky M., Toze S., 2017. Amplicon-based profiling of bacteria
587 in raw and secondary treated wastewater from treatment plants across Australia. Appl
588 Microbiol Biotechnol. 101(3), 1253-1266.
- 589 Anonymous, 2009. Ponds for stabilising organic matter.
590 www.water.wa.gov.au/__data/assets/pdf_file/0005/4100/84601.pdf
- 591 Apprill, A., McNally, S., Parsons, R., Weber, L., 2015. Minor revision to V4 region SSU rRNA
592 806R gene primer greatly increases detection of SAR11 bacterioplankton. Aquat. Microb.
593 Ecol. 75, 129-137.
- 594 Baars, J.K., 1962. The use of oxidation ditches for treatment of sewage for small communities.
595 Bull. World Health Organ. 26, 465-474.
- 596 Balvočiūtė, M., Huson, D.H., 2017. SILVA, RDP, Greengenes, NCBI and OTT—how do these
597 taxonomies compare? BMC Genomics. 18, 114.
- 598 Beilfuss, H.A., Quig, D., Block, M.A., Schreckenberger, P.C., 2015. Definitive identification of
599 *Laribacter hongkongensis* acquired in the United States. J. Clin. Microbiol. 53, 2385-2388.
- 600 Bellucci, M., Curtis, T.P., 2011. Ammonia-oxidizing bacteria in wastewater. Methods Enzymol.
601 496, 269-286.
- 602 Berlec, A., 2012. Novel techniques and findings in the study of plant microbiota: Search for plant
603 probiotics. Plant Sci. 193-194,96-102.
- 604 Bhaduri, S., Wesley, I.V., Bush, E.J., 2005. Prevalence of pathogenic *Yersinia enterocolitica*
605 strains in pigs in the United States. Appl Environ Microbiol. 71, 7117-7121.

606 Bokulich, N.A., Kaehler, B.D., Rideout, J.R., Dillon, M., Bolyen, E., Knight, R., Huttley, G.A.,
607 Caporaso, J.G., 2018. Optimizing taxonomic classification of marker-gene amplicon
608 sequences with QIIME 2's q2-feature-classifier plugin. *Microbiome*. 6, 90.

609 Brosius, J., Palmer, M. L., Kennedy, P. J., Noller, H. F., 1978. Complete nucleotide sequence of a
610 16S ribosomal RNA gene from *Escherichia coli*. *Proc Natl Acad Sci U S A*. 75, 4801-
611 4805.

612 Cabral, J.P., 2010. Water microbiology. Bacterial pathogens and water. *Int J Environ Res Public*
613 *Health*. 7, 3657-3703.

614 Caporaso, J.G., Kuczynski, J., Stombaugh, J., Bittinger, K., Bushman, F.D., Costello, E.K., Fierer,
615 N., Peña, A.G., Goodrich, J.K., Gordon, J.I., Huttley, G.A., Kelley, S.T., Knights, D.,
616 Koenig, J.E., Ley, R.E., Lozupone, C.A., McDonald, D., Muegge, B.D., Pirrung, M.,
617 Reeder, J., Sevinsky, J.R., Turnbaugh, P.J., Walters, W.A., Widmann, J., Yatsunenko, T.,
618 Zaneveld, J., Knight. R., 2010. QIIME allows analysis of high-throughput community
619 sequencing data. *Nat Methods*. 7, 335-336.

620 Caporaso, J. G., Lauber, C. L., Walters, W. A., Berg-Lyons, D., Lozupone, C. A., Turnbaugh, P.
621 J., Noah Fierer, N., Knight, R., 2011. Global patterns of 16S rRNA diversity at a depth of
622 millions of sequences per sample. *Proc. Natl. Acad. Sci. USA*. 108, 4516-4522.

623 Chakraborty, R. O'Connor, S.M. Chan, E. Coates, J.D., 2005. Anaerobic degradation of benzene,
624 toluene, ethylbenzene, and xylene compounds by *Dechloromonas* strain RCB *Appl.*
625 *Environ. Microbiol.* 71, 8649-8655.

626 Chakravorty, S., Helb, D., Burday, M., Connell, N., Alland, D., 2007. A detailed analysis of 16S
627 ribosomal RNA gene segments for the diagnosis of pathogenic bacteria. *J Microbiol*
628 *Methods*. 69, 330-339.

629 Connor, R., Renata, A., Ortigara, C., Koncagül, E., Uhlenbrook, S., Lamizana-Diallo, B.M.,
630 Zadeh, S.M., Qadir, M., Kjellén, M., Sjödin, J., Hendry, S., 2017. The United Nations
631 world water development report. Wastewater: the untapped resource. Paris, UNESCO.
632 <https://reliefweb.int/sites/reliefweb.int/files/resources/247153e.pdf>.

633 Cydzik-Kwiatkowska, A., Zielińska, M., 2016. Bacterial communities in full-scale wastewater
634 treatment systems. *World J Microbiol Biotechnol.* 32, 66.

635 Devanga Ragupathi, N.K., Muthuirulandi Sethuvel, D.P., Inbanathan, F.Y., Veeraraghavan, B.,
636 2017. Accurate differentiation of *Escherichia coli* and *Shigella serogroups*: challenges and
637 strategies. *New Microbes New Infect.* 21, 58-62.

638 Edgar, R.C., 2010. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics.*
639 26, 2460-2461.

640 Edgar, R.C., 2016. UNOISE2: improved error-correction for Illumina 16S and ITS amplicon
641 sequencing. *bioRxiv.* 081257.

642 Eland, L.E., Davenport, R.J., Santos, A.B., Mota Filho, C.R., 2018. Molecular evaluation of
643 microalgal communities in full-scale waste stabilisation ponds. *Environ. Technol.* In press.
644 doi: 10.1080/09593330.2018.1435730.

645 Fukushima, M., Kakinuma, K., Kawaguchi, R., 2002. Phylogenetic analysis of *Salmonella*,
646 *Shigella*, and *Escherichia coli* strains on the basis of the *gyrB* gene sequence. *J Clin*
647 *Microbiol.* 40, 2779-2785.

648 Garrido-Cardenas, J.A., Polo-López, M.I., Oller-Alberola, I., 2017. Advanced microbial analysis
649 for wastewater quality monitoring: metagenomics trend. *Appl. Microbiol. Biotechnol.* 101,
650 7445-7458.

651 Goodwin, S., McPherson, J.D., McCombie, W.R., 2106. Coming of age: ten years of next-
652 generation sequencing technologies. *Nat Rev Genet.* 17, 333-351.

653 Guentzel, M.N., 1996. *Escherichia, Klebsiella, Enterobacter, Serratia, Citrobacter, and Proteus.*
654 In: Baron S, editor. *Medical Microbiology.* 4th edition. Galveston (TX): University of
655 Texas Medical Branch at Galveston.

656 Guindon, S., Dufayard, J.F., Lefort, V., Anisimova, M., Hordijk, W., Gascuel, O., 2010. New
657 algorithms and methods to estimate maximum-likelihood phylogenies: assessing the
658 performance of PhyML 3.0. *Syst. Biol.* 59, 307-321.

659 Haft, D.H., Tovchigrechko, A., 2012. High-speed microbial community profiling. *Nat Methods.*
660 9, 793-794.

661 Hillmann, B., Al-Ghalith, G. A., Shields-Cutler, R.R., Zhu, Q., Gohl, D.M., Beckman, K.B.,
662 Knight, R., Knights, D., 2018. Evaluating the information content of shallow shotgun
663 metagenomics. *mSystems.* 3, e00069-00018.

664 Ho, L.T., Van Echelpoel, W., Goethals, P.L.M., 2017. Design of waste stabilization pond systems:
665 a review. *Water Res.* 123, 236-248.

666 Ibal, J.C., Pham, H.Q., Park, C.E., Shin, J-H., 2019. Information about variations in multiple copies
667 of bacterial 16S rRNA genes may aid in species identification. *PLoS One.* 14, e0212090.

668 Inaba, T., Hori, T., Aizawa, H., Ogata, A., Habe, H., 2017. Architecture, component, and
669 microbiome of biofilm involved in the fouling of membrane bioreactors. *NPJ Biofilms*
670 *Microbiomes.* 3, 5.

671 Katoh, K., Misawa, K., Kuma, K.-i., Miyata, T., 2002. MAFFT: a novel method for rapid multiple
672 sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* 30, 3059-3066.

673 Kayman, T., Abay, S., Hizlisoy, H., Atabay, H.İ., Diker, K.S., Aydin, F., 2012. Emerging pathogen
674 *Arcobacter* spp. in acute gastroenteritis: molecular identification, antibiotic susceptibilities
675 and genotyping of the isolated arcobacters. J Med Microbiol. 61, 1439-1444.

676 Kearse, M., Moir, R., Wilson, A., Stones-Havas, S., Cheung, M., Sturrock, S., Buxton, S., Cooper,
677 A., Markowitz, S., Duran, C., Thierer, T., Ashton, B., Meintjes, P., Drummond, A., 2012.
678 Geneious Basic: an integrated and extendable desktop software platform for the
679 organization and analysis of sequence data. Bioinformatics. 28, 1647-1649.

680 Kembel, S.W., Wu, M., Eisen, J.A. and Green, J.L., 2012. Incorporating 16S gene copy number
681 information improves estimates of microbial diversity and abundance. PLoS computational
682 biology, 8(10), p.e1002743.

683 Kim, Y., Koh, I., Rho, M., 2015. Deciphering the human microbiome using next-generation
684 sequencing data and bioinformatics approaches. Methods. 79–80, 52-59.

685 Kimura, M., 1980. A simple method for estimating evolutionary rate of base substitutions through
686 comparative studies of nucleotide sequences. J. Mol. Evol. 16, 111-120.

687 Levican, A., Collado, L., Aguilar, C., Yustes, C., Diéguez, A.L., Romalde, J.L. Figueras, M.J.,
688 2012. *Arcobacter bivalviorum* sp. nov. and *Arcobacter venerupis* sp. nov., new species
689 isolated from shellfish. Syst Appl Microbiol. 35, 133-138.

690 Lu, X., Zhang, X.-X., Wang, Z., Huang, K., Wang, Y., Liang, W., Tan, Y., Liu, B., Tang, J., 2015.
691 Bacterial pathogens and community composition in advanced sewage treatment systems
692 revealed by metagenomics analysis based on high-throughput sequencing. PLoS One. 10,
693 e0125549.

694 McDonald, D., Price, M.N., Goodrich, J., Nawrocki, E.P., DeSantis, T.Z., Probst, A., Andersen,
695 G.L., Knight, R., Hugenholtz, P., 2012. An improved Greengenes taxonomy with explicit
696 ranks for ecological and evolutionary analyses of bacteria and archaea. *ISME J.* 6, 610.

697 McLellan, S., Huse, S., Mueller Spitz, S., Andreishcheva, E., Sogin, M., 2010. Diversity and
698 population structure of sewage derived microorganisms in wastewater treatment plant
699 influent. *Environ. Microbiol.* 12, 378-392.

700 Oksanen, J., Blanchet, F.G., Friendly, M., Kindt, R., Legendre, P., McGlenn, D., Minchin, P.R.,
701 O'Hara, R. B., Simpson, G.L., Solymos, P., Henry, M., Stevens, H., Szoecs, E., Wagner, H.,
702 2018. *Vegan: Community ecology package*. R package version 2.5-2. [https://CRAN.R-](https://CRAN.R-project.org/package=vegan)
703 [project.org/package=vegan](https://CRAN.R-project.org/package=vegan).

704 Parada, A.E., Needham, D.M., Fuhrman, J.A., 2015. Every base matters: assessing small subunit
705 rRNA primers for marine microbiomes with mock communities, time series and global
706 field samples. *Environ. Microbiol.* 18, 1403-1414.

707 R Core Team, 2013. *R: A language and environment for statistical computing*. R Foundation for
708 Statistical Computing, Vienna, Austria. URL <http://www.R-project.org/>.

709 Rajasulochana, R., Preethy, V., 2016. Comparison on efficiency of various techniques in treatment
710 of waste and sewage water – a comprehensive review. *Resource-Efficient Technologies.* 2,
711 175-184.

712 Rodriguez-Manzano, J., Alonso, J.L., Ferrus, M.A., Moreno, Y., Amoros, I., Calgua, B., Hundesa,
713 A., Guerrero-Latorre, L., Carratala, A., Rusinol, M., Girones, R., 2012. Standard and new
714 faecal indicators and pathogens in sewage treatment plants, microbiological parameters for
715 improving the control of reclaimed water. *Water Sci. Technol.* 66, 2517-2523.

716 Safton, S., Cooper, G., Harrison, M., Wright, L. and Walsh, P., 1999. *Neisseria canis* infection: a
717 case report. *Commun Dis Intell.* 23, 221.

718 Sanapareddy, N., Hamp, T.J., Gonzalez, L.C., Hilger, H.A., Fodor, A.A., Clinton, S.M., 2009.
719 Molecular diversity of a North Carolina wastewater treatment plant as revealed by
720 pyrosequencing. *Appl Environ Microbiol.* 75, 1688-1696.

721 Seng, P., Boushab, B.M., Romain, F., Gouriet, F., Bruder, N., Martin, C., Paganelli, F., Bernit, E.,
722 Le Treut, Y.P., Thomas, P., Papazian, L., 2016. Emerging role of *Raoultella ornithinolytica*
723 in human infections: a series of cases and review of the literature. *Int J Infect Dis.* 45, 65-
724 71.

725 Sharafi, K., Davil, M.F., Heidari, M., Almasi, A., Taheri, H. 2012. Comparison of conventional
726 activated sludge system and stabilization pond in removal of chemical and biological
727 parameters. *Int. J. Environ. Health Eng.* 1, 1-5.

728 Shchegolkova, N.M., Krasnov, G.S., Belova, A.A., Dmitriev, A.A., Kharitonov, S.L., Klimina,
729 K.M., Melnikova, N.V., Kudryavtseva, A.V., 2016. Microbial community structure of
730 activated sludge in treatment plants with different wastewater compositions *Front*
731 *Microbiol.* 7, 90.

732 Shinoda, Y., Sakai, Y., Uenishi, H., Uchihashi, Y., Hiraishi, A., Yukawa, H., Yurimoto, H., Kato,
733 N., 2004. Aerobic and anaerobic toluene degradation by a newly isolated denitrifying
734 bacterium, *Thauera* sp. strain DNT-1 *Appl. Environ. Microbiol.* 70, 1385- 1392.

735 Siripong, S. and Rittmann, B.E., 2007. Diversity study of nitrifying bacteria in full-scale municipal
736 wastewater treatment plants. *Water Res.* 41, 1110-1120.

737 Stamatakis, A., 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
738 phylogenies. *Bioinformatics.* 30, 1312-1313.

739 Verbyla, M., von Sperling, M., Maiga, Y., 2017. Waste Stabilization Ponds. In: J.B. Rose and B.
740 Jiménez- Cisneros, (eds) Global Water Pathogens Project. <http://www.waterpathogens.org>
741 (C. Haas, J.R. Mihelcic and M.E. Verbyla) (eds) Part 4 Management of risk from Excreta
742 and Wastewater) www.waterpathogens.org/book/waste-stabilization-ponds Michigan
743 State University, E. Lansing, MI, UNESCO.

744 Von Sperling, M., 2007. Waste stabilisation ponds. IWA Publishing.
745 www.iwapublishing.com/sites/default/files/ebooks/9781780402109.pdf.

746 Xu, S., Yao, J., Ainiwaer, M., Hong, Y., Zhang, Y., 2018. Analysis of bacterial community
747 structure of activated sludge from wastewater treatment plants in winter. *Biomed Res Int.*
748 2018, 8278970.

749 Yang, C., Zhang, W, Liu, R, Li, Q, Li, B, Wang, S, Song, C, Qiao, C, Mulchandani, A., 2011.
750 Phylogenetic diversity and metabolic potential of activated sludge microbial communities
751 in full-scale wastewater treatment plants. *Environ Sci Technol.* 45, 7408-7415.

752 Ye, L., Zhang, T., 2011. Pathogenic bacteria in sewage treatment plants as revealed by 454
753 pyrosequencing. *Environ Sci Technol* 45, 7173-7179.

754 Zhang, T., Shao, M.-F., Ye, L., 2012. 454 pyrosequencing reveals bacterial diversity of activated
755 sludge from 14 sewage treatment plants. *ISME J.* 6, 1137-1147.

756

757

758 **Figure Legends**

759 **Figure 1.** WWTP localities and different treatment stages sampled.

760 **Figure 2.** 16S NGS sequence percent composition plot of phyla (P) and classes (C) detected in
761 different treatment stages of wastewater sampled from WWTPs 1-4. Treatment stages include
762 influent (I), intermediate (INT) and effluent (E). Phyla with $\leq 10\%$ overall sequence composition
763 are grouped as “other”.

764 **Figure 3.** Maximum likelihood tree of a 247 bp alignment (gaps excluded) of genomic 16S
765 Enterobacteriaceae sequences trimmed to the V4 region. The seven Enterobacteriaceae ZOTU
766 sequences derived from this study are in bold typeface. Values at nodes indicate Bootstrap values
767 from 1,000 replicates. Outgroup of tree *Vibrio cholerae* (2614873) not shown.

768

769 **Appendices**

770 **Appendix A. Supplementary Figures**

771 **Figure A.1.** Alpha rarefaction plot of 16S sequencing depth and ZOTUs detected in WWTP
772 samples prior to low read abundance (<0.05%) filtering.

773 **Figure A.2.** Alpha rarefaction plots of 16S sequencing depth and ZOTUs detected prior to low
774 read abundance (<0.05%) filtering for WWTPs 1-4 and treatment stages.

775 **Figure A.3.** Alpha rarefaction plot of 16S sequencing depth and ZOTUs detected in WWTP
776 samples after low read abundance (<0.05%) filtering.

777 **Figure A.4.** Alpha rarefaction plots of 16S sequencing depth and ZOTUs detected after low read
778 abundance (<0.05%) filtering for WWTPs 1-4 and treatment stages.

779

780 **Appendix B. Supplementary Data**

781 **Supplementary File B.1.** List of 3,598 16S V4 region ZOTU sequences generated by this study.

782 **Supplementary File B.2.** Sequence totals and compositions.

783 **Supplementary File B.3.** Pairwise genetic distance matrix of the 247 bp alignment (gaps
784 excluded) of genomic 16S Enterobacteriaceae sequences trimmed to the V4 region that was used
785 to construct the phylogenetic tree in Figure 3.

786 **Supplementary File B.4.** Comparison of Greengenes and NCBI nr/nt database taxa to ZOTUs
787 potentially from pathogenic bacteria.

788 **Supplementary File B.5.** Comparison of Greengenes and NCBI nr/nt database taxa to ZOTUs
789 from nitrifying, denitrifying and floc-forming bacteria.

790 **Supplementary File B.6.** Chloroplast sequences in the Greengenes 99 OTU taxonomy file
791 renamed with the kingdom “Chloroplast”.

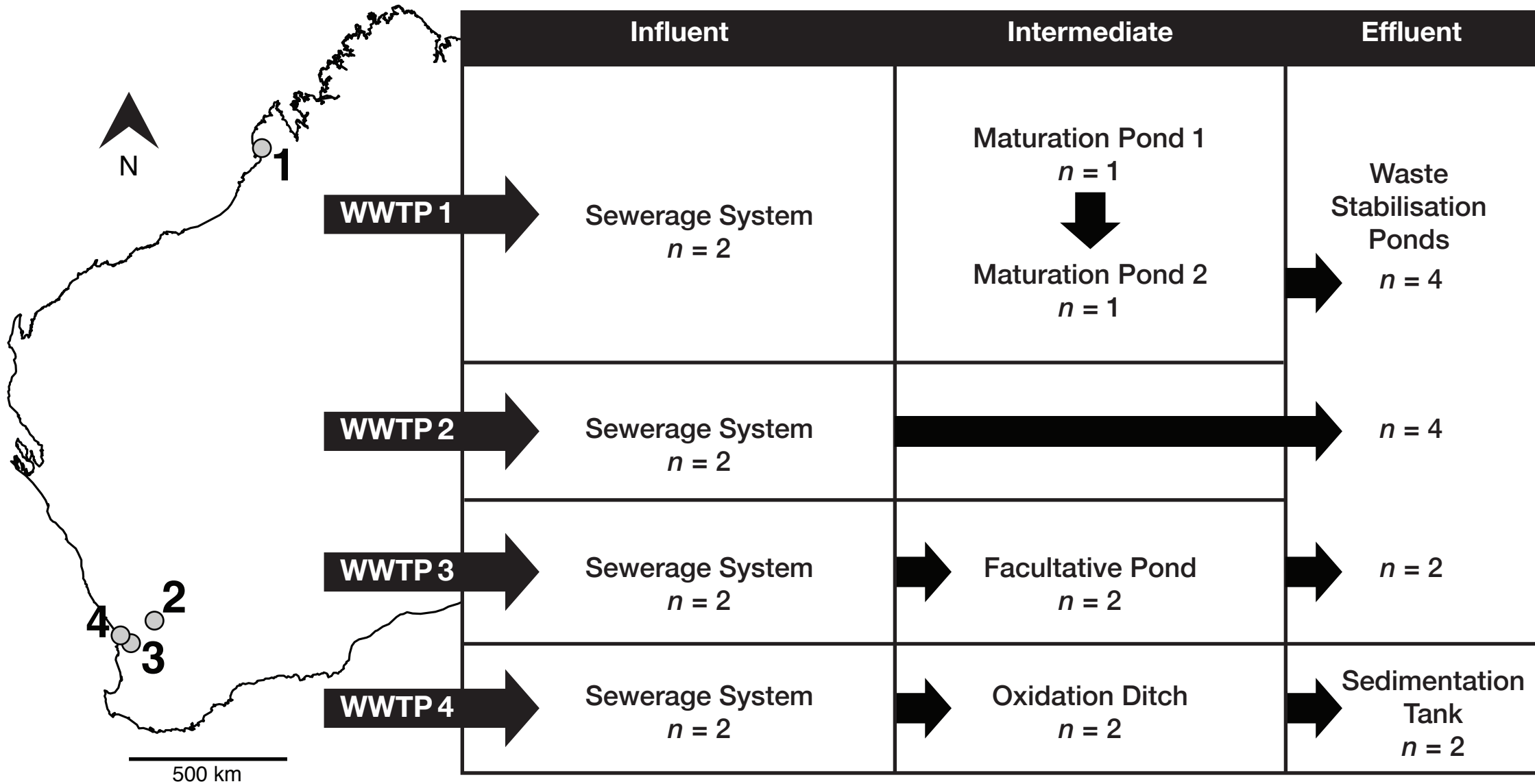
792

793

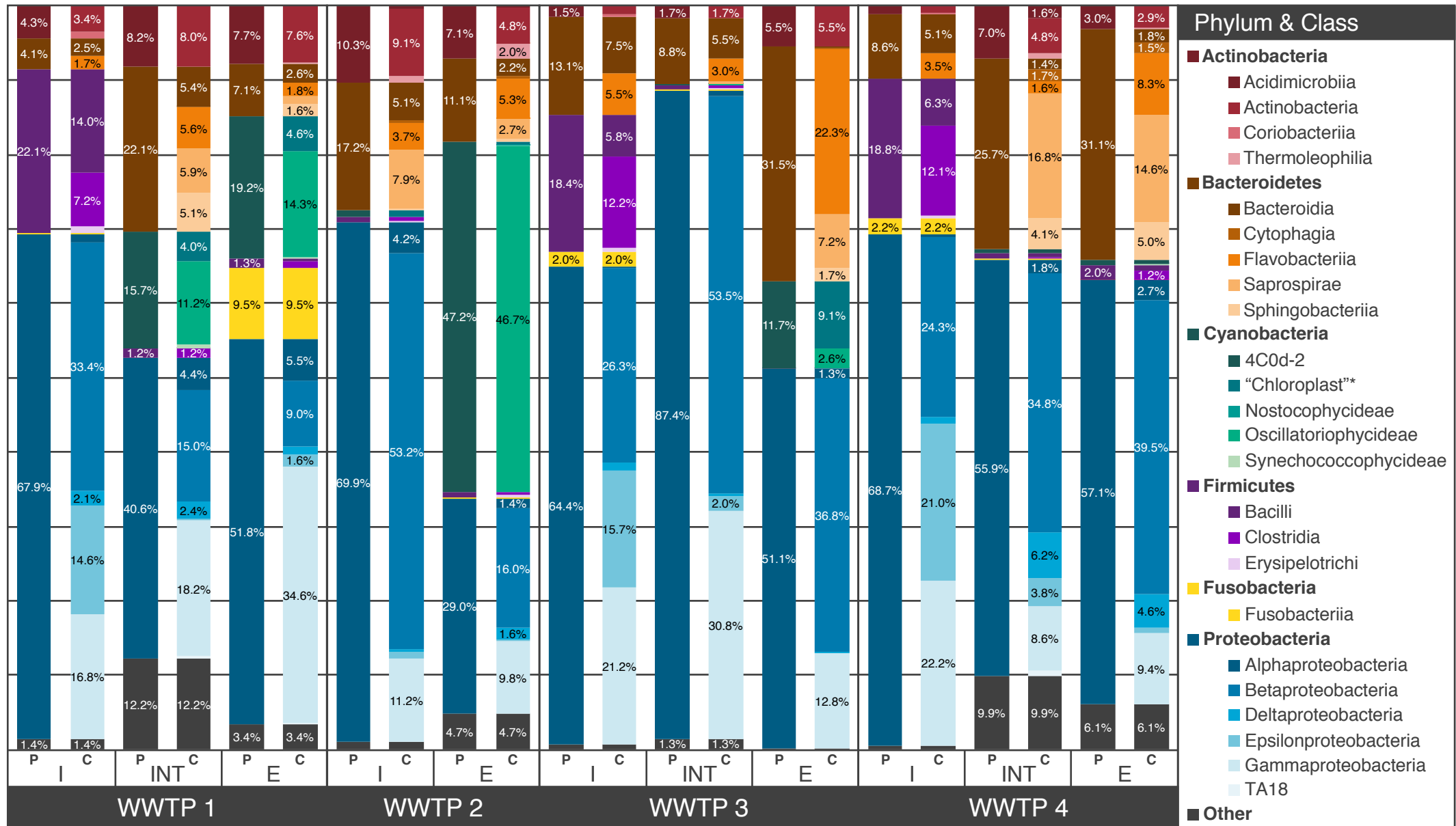
794

795

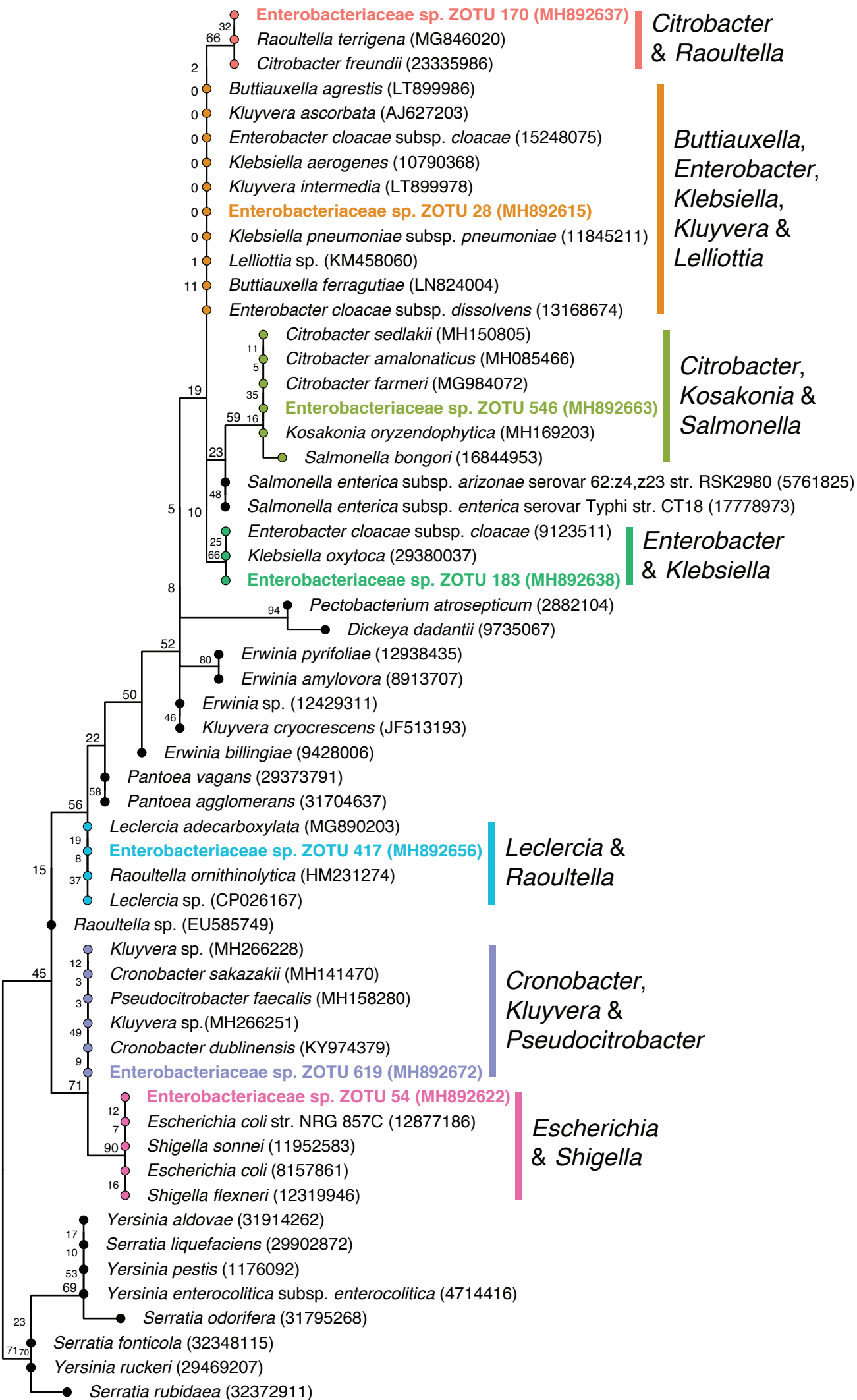
796



Percent composition of 16S sequences (%)



Wastewater treatment plant and treatment stage



0.02

Table 1. Rural wastewater treatment plant samples analysed in the present study.

WWTP	Treatment technology	Location	Climate	Sample ID	Wastewater treatment stage	Sample collection date; season		
WWTP 1	Stabilisation pond: Combined anaerobic and aerobic pond system, followed by two maturation ponds	Northwest Western Australia	Tropical climate. Wet and dry seasons.	WWTP 1-1	Influent	19-Feb-2015; Wet		
				WWTP 1-2	Effluent (pre-chlorination)			
				WWTP 1-3	Effluent (post-chlorination)			
				WWTP 1-4	Influent	7-Sep-2015; Dry		
							WWTP 1-5	Intermediate (post maturation pond 1)
							WWTP 1-6	Intermediate (post maturation pond 2)
							WWTP 1-7	Effluent (pre-chlorination)
							WWTP 1-8	Effluent (post-chlorination)
WWTP 2	Stabilisation pond: One facultative pond	Wheatbelt, Western Australia	Hot dry summers and mild winters. Four distinct seasons.	WWTP 2-1	Influent	12-Feb-2015; Summer		
				WWTP 2-2	Effluent (final pond)			
				WWTP 2-3	Effluent (storage basin)			
				WWTP 2-4	Influent	13-Jul-2015; Winter		
							WWTP 2-5	Effluent (final pond)
							WWTP 2-6	Effluent (storage basin)

WWTP 3	Stabilisation pond: Two primary facultative ponds, and one secondary pond	Southwest Western Australia	Temperate climate. Four distinct seasons.	WWTP 3-1	Influent	23-Feb-2015; Summer
				WWTP 3-2	Intermediate (post-pond)	
				WWTP 3-3	Effluent	
				WWTP 3-4	Influent	14-July-2015; Winter
				WWTP 3-5	Intermediate (post-pond)	
				WWTP 3-6	Effluent	
WWTP 4	Activated sludge: Oxidation ditches followed by sedimentation tanks	Southwest Western Australia	Temperate climate. Four distinct seasons.	WWTP 4-1	Influent	23-Feb-2015; Summer
				WWTP 4-2	Intermediate (oxidation ditch)	
				WWTP 4-3	Effluent	
				WWTP 4-4	Influent	14-July-2015; Winter
				WWTP 4-5	Intermediate (oxidation ditch)	
				WWTP 4-6	Effluent	

Table 2. V4 16S NGS sequence statistics.

Statistics	Raw (unprocessed)	Pre-processed ^a	Processed 16S sequences ^b			
	Grand total (<i>n</i> = 34)		Samples (<i>n</i> = 26)	Extraction controls (<i>n</i> = 6)	NTCs (<i>n</i> = 2)	Grand total (<i>n</i> = 34)
Average	27,965	23,805	26,746	8	8	20,454
Standard deviation	27,254	24,239	20,608	7	2	21,314
Min	2,646	2	4,681	2	6	2
Max	182,113	95,135	85,305	21	9	85,305
Total	1,426,191	809,368	695,400	48	19	695,463

^aMerged, quality filtered sequences with singletons and chimeras removed

^bMerged, quality filtered sequences with singletons, chimeras, unassigned sequences and low abundance sequences (<0.05%) removed

Table 3. Enterobacteriaceae (Gammaproteobacteria: Enterobacteriales) ZOTUs Greengenes assigned taxonomy cross-checked against the NCBI nr/nt database.

ZOTU no.	Accession no.	Final taxonomy	Greengenes results		NCBI nr/nt results			Correct Greengenes taxonomy?		
			Assigned taxonomy	Confidence scores ^a	GenBank® accession no.	Species	Percent identity	Family	Genus	Species
28	MH892615	Enterobacteriaceae sp.	Enterobacteriaceae sp.	0.95	MH384426	<i>Enterobacter xiangfangensis</i>	100	✓	*	*
					MH190220	<i>Erwinia aphidicola</i>	100	✓	*	*
					MH411220	<i>Klebsiella pneumoniae</i>	100	✓	*	*
54	MH892622	Enterobacteriaceae sp.	<i>Escherichia coli</i>	0.96	MH396737	<i>Escherichia coli</i>	100	✓	×	×
					MH352164	<i>Salmonella enterica</i> subsp. <i>enterica</i>	100	✓	×	×
					MH371327	<i>Shigella flexneri</i>	100	✓	×	×
170	MH892637	Enterobacteriaceae sp.	<i>Citrobacter</i> sp.	0.95	NR156052	<i>Citrobacter europaeus</i>	100	✓	×	*
					MH371322	<i>Citrobacter freundii</i>	100	✓	×	*
					MH352205	<i>Salmonella enterica</i> subsp. <i>enterica</i>	100	✓	×	*
183	MH892638	Enterobacteriaceae sp.	Enterobacteriaceae sp.	0.94	CP020089	<i>Enterobacter cloacae</i>	100	✓	*	*
					MF360016	<i>Klebsiella michiganensis</i>	100	✓	*	*
					MH196342	<i>Klebsiella oxytoca</i>	100	✓	*	*
417	MH892656	Enterobacteriaceae sp.	Enterobacteriaceae sp.	0.92	MG890203	<i>Leclercia adecarboxylata</i>	100	✓	*	*
					MG022656	<i>Raoultella electrica</i>	100	✓	*	*
					MG516115	<i>Raoultella ornithinolytica</i>	100	✓	*	*
546	MH892663	Enterobacteriaceae sp.	<i>Trabulsiella</i> sp.	1	MH085457	<i>Citrobacter amalonaticus</i>	100	✓	×	*
					MF186607	<i>Citrobacter farmeri</i>	100	✓	×	*
					MH169203	<i>Kosakonia oryzendophytica</i>	100	✓	×	*
619	MH892672	Enterobacteriaceae sp.	Enterobacteriaceae sp.	1	MH141470	<i>Cronobacter sakazakii</i>	100	✓	*	*
					MH169205	<i>Kluyvera georgiana</i>	100	✓	*	*
					MG890202	<i>Pseudocitrobacter faecalis</i>	100	✓	*	*

*Taxon was unassigned, which was the correct choice based on BLAST results.

^aConfidence scores are probabilities generated by the naïve Bayes algorithm implemented by QIIME 2 feature classifier (<https://docs.qiime2.org/2018.6/tutorials/feature-classifier/>).

Table 4. Sequence composition (%) of pathogens and possible pathogens in WWTPs 1-4 influent (I), intermediate (INT) and effluent (E) with taxonomy confirmed with Greengenes and NCBI nr/nt sequence databases.

Class	Order	Family	Taxonomic assignment ^a	ZOTU no.	Accession no.	WWTP 1			WWTP 2		WWTP 3			WWTP 4		
						I	INT	E	I	E	I	INT	E	I	INT	E
Actinobacteria																
Actinobacteria	Actinomycetales	Corynebacteriaceae	<i>Corynebacterium</i> sp.	2603	MH892704	<0.1	-	-	-	-	-	-	-	-	-	-
		Mycobacteriaceae	<i>Mycobacterium</i> spp.	36; 332; 588; 1469; 1651; 1756; 1801	MH892617; MH892651; MH892668; MH892692; MH892696; MH892698; MH892699	-	2.5	2.2	0.1	-	-	0.1	0.3	-	0.2	0.1
Chlamydiae																
Chlamydiia	Chlamydiales	Parachlamydiaceae	Parachlamydiaceae sp.	1459	MH892691	-	-	-	-	-	-	-	-	-	0.1	0.1
		Rhabdochlamydiaceae	<i>Candidatus Rhabdochlamydia</i> sp.	3044	MH892708	-	-	-	-	-	-	-	-	-	-	<0.1
		-	Chlamydiales spp.	1597; 2741; 3295; 3540	MH892695; MH892705; MH892711; MH892712	-	-	-	-	-	-	-	-	-	-	0.2
	-	Chlamydiia spp.	3035; 3270	MH892707; MH892710	-	-	-	-	-	-	-	-	-	-	-	0.1
Firmicutes																
Bacilli	Lactobacillales	Streptococcaceae	Lactobacillales spp.	31; 134; 443	MH892616; MH892632; MH892657	1.5	-	-	-	-	2.3	-	-	1.9	0.2	0.2
			Streptococcaceae sp.	910	MH892678	0.1	-	-	-	-	-	-	-	-	-	-
			<i>Streptococcus</i> spp.	8; 559	MH892611; MH892665	12.1	-	0.2	-	-	2.6	-	-	3.2	0.4	0.2
Clostridia	Clostridiales	Clostridiaceae	Clostridiaceae spp.	357; 1161	MH892652; MH892684	-	-	0.2	0.1	-	-	-	-	-	-	
			Ruminococcaceae sp.	1387	MH892688	<0.1	-	-	-	-	-	-	-	-	-	
		-	Clostridiales spp.	359; 460; 1119; 1967	MH892653; MH892659; MH892683; MH892701	0.1	0.4	-	-	-	0.1	-	-	0.1	-	-
Bacilli	Lactobacillales	Enterococcaceae	<i>Enterococcus</i> spp.	286; 288; 1246	MH892646; MH892648; MH892686	0.2	-	0.1	-	-	0.2	-	-	0.3	-	-
Proteobacteria																
Betaproteobacteria	Burkholderiales	Alcaligenaceae	Alcaligenaceae sp.	39	MH892619	-	-	<0.1	0.2	2.0	-	-	-	-	-	-
	Neisseriales	Neisseriaceae	<i>Laribacter hongkongensis</i>	55	MH892623	0.8	-	-	-	-	0.9	0.1	-	0.7	0.1	0.1
			<i>Neisseria canis</i>	74	MH892626	1.1	-	-	-	-	0.6	-	-	0.3	<0.1	0.1
			<i>Neisseria</i> sp.	1251	MH892687	<0.1	-	-	-	-	-	-	-	-	-	-
			Neisseriaceae spp.	19; 197; 287;	MH892614; MH892639;	0.9	-	-	0.1	-	4.2	0.3	-	4.0	0.2	0.1

				960	MH892647; MH892679														
			<i>Vitreoscilla</i> spp.	97; 169; 466	MH892629; MH892636; MH892660	0.2	-	-	-	-	1.0	0.1	-	1.0	-	-			
Epsilonproteobacteria	Campylobacterales	Campylobacteraceae	<i>Arcobacter</i> spp.	218; 289; 598; 1174	MH892642; MH892649; MH892670; MH892685	0.3	-	-	-	-	0.6	-	-	0.1	-	-			
			<i>Arcobacter venerupis</i>	214	MH892641	0.7	-	-	-	-	-	-	-	-	-	-	-	-	
			Campylobacteraceae spp.	2; 37; 59; 158; 229	MH892609; MH892618; MH892624; MH892635; MH892643	13.4	0.2	1.6	0.9	0.1	14.9	1.9	-	20.7	3.8	0.7			
Gammaproteobacteria	Aeromonadales	Aeromonadaceae	Aeromonadaceae spp.	483; 639; 1476	MH892661; MH892673; MH892693	0.2	-	-	-	-	0.1	-	-	-	-	-			
			<i>Aeromonas</i> sp.	3	MH892610	6.0	2.2	8.6	0.4	0.2	4.7	1.1	4.8	2.8	0.5	0.8			
	Enterobacteriales	Enterobacteriaceae	Enterobacteriaceae spp.	28; 54; 170; 183; 417; 546; 619	MH892615; MH892622; MH892637; MH892638; MH892656; MH892663; MH892672	3.7	0.2	0.6	-	-	2.6	-	0.6	2.4	0.1	0.3			
	Legionellales	Coxiellaceae	Coxiellaceae sp.	892	MH892677	-	-	-	-	-	-	-	-	-	-	-	-	0.4	
		-	Legionellales sp.	3079	MH892677	-	-	-	-	<0.1	-	-	-	-	-	-	-	-	
		Legionellaceae	<i>Legionella</i> sp.	2554	MH892703	-	-	-	-	-	-	-	-	-	-	-	-	0.1	
	Pseudomonadales	Moraxellaceae	<i>Acinetobacter</i> spp.	10; 16; 41; 45; 75; 101; 283; 317; 564; 584; 886; 991; 992	MH892612; MH892613; MH892620; MH892621; MH892627; MH892630; MH892645; MH892650; MH892666; MH892667; MH892676; MH892681; MH892682	1.0	-	0.1	-	-	8.5	17.1	-	12.8	0.7	0.7			
				Pseudomonadaceae	<i>Pseudomonas</i> spp.	65; 77; 109; 147; 151; 210; 233; 361; 402; 515; 556; 607; 678; 722; 1402	MH892625; MH892628; MH892631; MH892633; MH892634; MH892640; MH892644; MH892654; MH892655; MH892662; MH892664; MH892671; MH892674; MH892675; MH892689	1.1	-	1.3	0.1	1.3	0.3	3.1	1.4	0.6	-	0.1	
		-	-	Gammaproteobacteria spp.	1440; 1709	MH892690; MH892697	-	-	-	-	-	-	-	-	-	-	0.2	0.1	
		Vibrionales	Pseudoalteromonadaceae	<i>Vibrio</i> sp.	589	MH892669	-	-	-	-	0.1	-	-	-	-	-	-		
Spirochaetes																			
Leptospirae	Leptospirales	Leptospiraceae	Leptospiraceae sp.	2146	MH892702	-	-	-	-	-	-	-	-	-	-	-	0.1		
			Spirochaetes sp.	2988	MH892706	-	-	-	-	-	-	-	-	-	-	-	<0.	-	

																1	
			<i>Turneriella</i> sp.	1946	MH892700	-	-	-	-	-	-	-	-	-	-	0.1	0.1
Spirochaetia	Spirochaetales	Spirochaetaceae	Spirochaetaceae spp.	445; 965; 1564	MH892658; MH892680; MH892694	-	0.1	0.1	0.1	-	-	-	-	-	-	-	-

^aMost specific level of taxonomy designated after comparing ZOTUs to Greengenes and NCBI nr/nt databases.

Table 5. Sequence composition (%) of nitrifying, denitrifying and floc-forming bacteria in WWTPs 1-4 influent (I), intermediate (INT) and effluent (E) with taxonomy confirmed with Greengenes and NCBI nr/nt sequence databases.

Class	Order	Family	ZOTU no.	Accession no.	Species	WWTP 1			WWTP 2		WWTP 3			WWTP 4			
						I	INT	E	I	E	I	INT	E	I	INT	E	
Bacteroidetes																	
Flavobacteriia	Flavobacteriales	Flavobacteriaceae	35; 44; 57; 103; 153; 155; 172; 178; 263; 365; 474; 660; 766; 970; 1142; 1241; 1883; 1912; 2042; 2242; 2349; 2374; 2375; 2493; 2649; 2905; 3231; 3371	MH892717; MH892718; MH892720; MH892728; MH892733; MH892734; MH892737; MH892738; MH892745; MH892752; MH892763; MH892771; MH892779; MH892792; MH892798; MH892800; MH892811; MH892813; MH892815; MH892816; MH892817; MH892818; MH892819; MH892820; MH892821; MH892823; MH892825; MH892826	<i>Flavobacterium</i> spp.	-	0.8	0.4	2.8	4.7	-	0.2	7.6	0.1	0.4	6.5	
Nitrospirae																	
Nitrospira	Nitrospirales	Nitrospira	404; 614; 1574; 2690	MH892758; MH892767; MH892808; MH892822	<i>Nitrospira</i> spp.	-	-	-	-	-	-	-	-	-	1.2	1.5	
Proteobacteria																	
Betaproteobacteria	Burkholderiales	Comamonadaceae	27	MH892715	<i>Aquabacterium</i> sp.	-	0.3	0.5	1.3	0.2	0.1	3.8	0.1	0.1	0.1	0.1	
			647	MH892769	<i>Brachymonas denitrificans</i>	0.1	-	-	-	-	-	-	-	-	-	-	-
			94; 677; 356; 492; 776; 580; 926; 1619; 3101; 855	MH892726; MH892772; MH892751; MH892765; MH892780; MH892766; MH892790; MH892809; MH892824; MH892785	Comamonadaceae spp.	0.3	0.1	0.1	-	0.3	0.6	0.6	-	0.4	0.6	1.1	
			11	MH892713	Comamonas sp.	6.9	-	0.1	0.1	-	4.4	1.4	-	4.4	0.2	0.2	
			1336	MH892802	Delftia sp.	-	-	-	-	-	-	-	-	-	-	-	
			73; 85; 100; 184; 225; 275; 319; 426; 449; 454; 1420	MH892723; MH892724; MH892727; MH892741; MH892744; MH892747; MH892749; MH892759; MH892761; MH892762; MH892805	<i>Hydrogenophaga</i> spp.	-	0.2	0.6	0.4	2.4	-	2.8	0.3	<0.1	-	<0.1	
			1109; 1403	MH892796; MH892804	<i>Polaromonas</i>	-	-	-	0.1	<0.	-	<0.1	-	-	-	0.2	

					spp.					1							
			716; 904	MH892774; MH892788	<i>Rhodoferax</i> spp.	-	0.2	-	-	-	-	-	-	-	-	0.5	
			762; 3498	MH892778; MH892828	<i>Rubrivivax</i> spp.	-	-	<0.1	-	0.1	-	-	-	-	-	-	
Rhodocyclales	Rhodocyclaceae		171; 823; 1898	MH892736; MH892782; MH892812	<i>Azoarcus</i> spp.	-	-	-	-	0.6	-	-	-	-	-	-	
			387; 980	MH892755; MH892794	<i>Azonexus</i> spp.	<0.1	-	-	-	0.2	-	-	-	-	-	-	-
			193; 863	MH892743; MH892786	<i>Azospira</i> spp.	0.1	<0.1	-	-	-	-	-	-	-	-	2.3	1.5
			11; 30; 71; 302; 490	MH892713; MH892716; MH892722; MH892748; MH892764	<i>Dechloromonas</i> spp.	1.8	-	0.1	0.1	-	0.4	5.6	-	0.3	1.0	1.2	
			1462	MH892806	<i>Methyloversatilis</i> sp.	-	-	-	-	<0.1	-	-	-	-	-	-	-
			49; 70; 104; 389; 717; 847; 901	MH892719; MH892721; MH892729; MH892756; MH892775; MH892784; MH892787	<i>Propionivibrio</i> spp.	2.2	-	-	0.2	-	2.1	0.9	-	1.8	0.3	0.1	
			623; 1139; 1305; 1367; 1546; 148; 269; 373; 754; 1151	MH892768; MH892797; MH892801; MH892803; MH892807; MH892732; MH892746; MH892753; MH892777; MH892799	Rhodocyclaceae spp.	0.6	<0.1	-	0.3	0.1	0.4	0.1	-	0.6	0.3	0.4	
			787; 385; 977	MH892781; MH892754; MH892793	<i>Sterolibacterium</i> spp.	-	-	-	-	-	-	-	-	-	-	1.3	1.2
			1796	MH892810	<i>Sulfuritalea</i> sp.	-	<0.1	-	-	-	-	-	-	-	-	-	-
			23; 91; 126; 924	MH892714; MH892725; MH892731; MH892789	<i>Thauera</i> spp.	3.4	-	<0.1	-	0.2	0.5	6.0	-	0.4	-	0.2	
			191; 1974	MH892742; MH892814	<i>Uliginosibacterium</i> spp.	-	-	<0.1	-	-	-	-	-	-	-	2.4	1.4
			120; 180; 181; 335	MH892730; MH892739; MH892740; MH892750	<i>Zoogloea</i> spp.	0.7	-	<0.1	0.1	-	0.8	0.6	-	0.5	0.1	1.0	
		Unclassified	Unclassified	165; 394; 441; 655; 693; 728; 825; 938; 1015; 3404	MH892735; MH892757; MH892760; MH892770; MH892773; MH892776; MH892783; MH892791; MH892795; MH892827	<i>Candidatus Accumulibacter</i> spp. ^b	-	-	-	-	-	-	-	-	-	5.0	5.9

^aMost specific level of taxonomy designated after comparing ZOTUs to Greengenes and NCBI nr/nt sequences.

^b*Candidatus Accumulibacter* spp. was assigned by Greengenes to the family Rhodocyclaceae, but is a recently discovered bacterium that has not yet been classified to an order or family.

