

RSC Advances

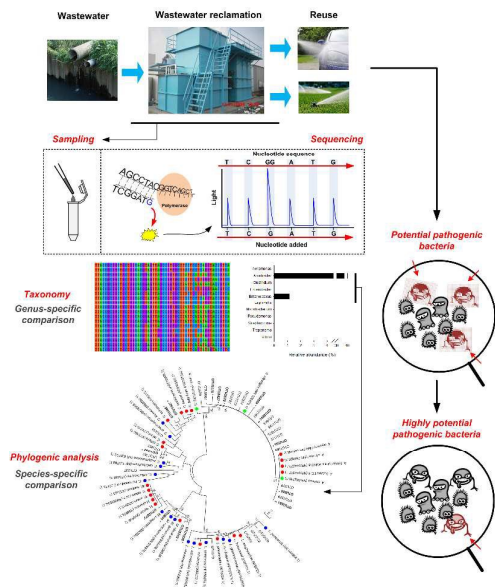


This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. This *Accepted Manuscript* will be replaced by the edited, formatted and paginated article as soon as this is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.



Efficient diagnosis based on pyrosequencing can address the highly potential pathogenic bacteria in wastewater.

1 **Occurrence and fate of potential pathogenic bacteria as revealed by**
2 **pyrosequencing in a full-scale membrane bioreactor treating**
3 **restaurant wastewater**

4 Jinxing Ma¹, Zhiwei Wang^{*1}, Lili Zang², Jian Huang¹, Zhichao Wu¹

5 ¹ State Key Laboratory of Pollution Control and Resource Reuse, School of Environmental
6 Science and Engineering, Tongji University, 1239 Siping Road, Shanghai 200092, P.R. China

7 ² Shanghai Zizheng Environm Technol Co Ltd, Shanghai 200437, P.R. China

8

9

10

11

12

13 Revised Manuscript Submitted to RSC Advances

14

15

16

17

18

19

20

21 *Corresponding Author. Tel./fax: +86-21-65980400 (Z. Wang); E-mail address:

22 zwwang@tongji.edu.cn (Z. Wang)

23 Abstract

24 One of the primary concerns on wastewater reuse is the presence of pathogenic bacteria.
25 Considering that indicator bacteria might only offer limited information, we applied
26 high-throughput pyrosequencing in this study to reveal bacterial pathogen diversity in a full-scale
27 membrane bioreactor (MBR) treating restaurant wastewater. The results showed that fecal
28 indicator bacteria could provide a rough estimation rather than an accurate characterization of the
29 potential pathogenic bacteria in wastewaters particularly from non-fecal sources. In general, MBR
30 treatment had a good removal of potential pathogenic bacteria. The bacterial counts of *Arcobacter*
31 was decreased by nearly seven orders of magnitude, from $(8.35 \pm 0.87) \times 10^7$ to <10 counts/mL,
32 and *Aeromonas*, *Enterobacter*, *Enterococcus*, and *Pseudomonas* were not detected in the treated
33 wastewater. The most dominant potential pathogens in activated sludge and treated wastewater
34 were affiliated to the genera of *Legionella*, *Clostridium* and *Mycobacterium*. Species-specific
35 comparison showed that only a small portion (0.0~1.6%) of the corresponding sequences had
36 identities of $> 99\%$ to the neighbor pathogenic species, including *Arcobacter butzleri* and
37 *Arcobacter cryaerophilus*. This study, therefore, provides insights into the occurrence and fate of
38 potential bacterial pathogens in restaurant wastewater treatment and reclamation using MBRs.

39 Keywords: pyrosequencing; membrane bioreactor; pathogenic bacteria; restaurant wastewater

40

41 1. Introduction

42 During the last decade, catering industry has experienced an explosive growth in China and
43 the business turnover was doubled from 2006 (\approx \$ 161 billion per year) to 2011 (\approx \$ 332 billion
44 per year)¹. Wastewater streams discharged from restaurants are generally characterized by high
45 content of oil, grease (O&G), suspended solid (SS) and detergent^{2,3}. High O&G is a tremendous
46 burden to the municipal wastewater systems because these organic substances usually tend to
47 clump together, causing drainage pipelines corrosion under anaerobic conditions. Therefore,
48 appropriate treatment of the restaurant wastewater is necessary in order to reduce the adverse
49 impacts of discharging⁴. As an option, the membrane bioreactor (MBR) is a fascinating and
50 promising technology, which presents distinctive advantages such as high volumetric organic
51 loading, small environmental footprint, and sound separation of emulsions that contain oil droplets
52 with diameter less than $20 \mu\text{m}$ ⁵. MBRs also offer the opportunity to spare the expenditure of

53 wastewater treatment since their superior effluent is more suitable for on-site reuse (e.g., flushing
54 toilets) in the restaurants.

55 On condition that treated wastewater utilization is expected, contaminant removals should be
56 sufficient to meet stringent regulatory standards, because of the public health concerns ⁶. In MBRs,
57 bacteria play an important role in the biochemical process, consuming nutrients and organic
58 matters. To date, numerous studies have been conducted to improve the efficiency of the
59 biochemical process, but most of them failed to attach importance to the potential hazard of the
60 bacteria accordingly. One of the primary concerns on bacteria is the community of pathogenic
61 bacteria originated from the excrement of disease-carrying humans and animals or other sources ⁷.
62 In many public places (e.g., general merchandise stores and restaurants), the outbreak of
63 gastroenteritis or other infections due to access of reused water could actually be masked by the
64 background levels of assumed sources, such as food-borne and community-based infections ⁸.
65 Microbial assessment of pathogenic bacteria in treated wastewater is thereby important in view of
66 consequent health risks.

67 Historically, fecal indicator bacteria including total and fecal coliforms and enterococci have
68 been widely used as a monitoring tool to predict the presence of potential bacterial, viral and
69 protozoan pathogens ⁹. The major drawback of fecal indicator bacteria arises from their poor
70 correlation with pathogens, especially those from non-fecal sources ^{8, 10}. Moreover, membranes
71 have a size-selective retention of different bacteria, and the abundance of pathogenic bacteria in
72 the permeate could be underestimated or overestimated when referred to certain indicator bacteria
73 ¹¹. In recent years, real-time qPCR assays have been proposed and these assays are now used in
74 many diagnostic and reference laboratories for the detection of pathogenic bacteria in clinical
75 fluids ^{8, 12, 13}. Compared to indicator bacteria methods, the qPCR assay enables quantitative and
76 highly specific detection, which could target 16S ribosomal RNA, encoding genes or
77 housekeeping genes of actual pathogens ¹². Nevertheless, the application of this technology is still
78 hindered due to its limited throughput capacity. In environmental samples, bacterial pathogen
79 diversity can be extremely high, as reflected by more than thirty phylogenetic genera and
80 thousands of strains. Clueless one-by-one detection is definitely time-consuming, which might
81 also miss the potential infectious risk. Therefore, illuminating solutions are urgently required for

82 elucidating bacterial pathogen diversity in MBRs and assessing full microbial risk of treated
83 wastewater reuse in public places.

84 In this study, 454 high-throughput pyrosequencing was used to investigate the occurrence and
85 fate of potential pathogenic bacteria in a full-scale MBR treating restaurant wastewater.
86 Pyrosequencing is a high-throughput analytical method that generates a large amount of DNA
87 reads through a massively parallel sequencing-by-synthesis approach, and this technology can
88 provide an adequate resolution to the microbial diversity of different environmental samples ^{7, 14-16}.
89 In the present work, 258,438 reads of the hypervariable V1~V3 regions of the bacterial 16S rRNA
90 gene were obtained. Sequence subsets with the capacity of 10,000 and 100 were generated from
91 the maternal gene libraries by a semi-random extraction method, and comprehensive comparison
92 of these datasets was then carried out. Bacterial pathogen diversity was analyzed at genus level
93 using Ribosomal Database Project (RDP) Classifier ¹⁷. Alignment of the corresponding sequences
94 to the known pathogen was further conducted by phylogenetic analysis.

95 **2. Materials and methods**

96 2.1. Sample collection and pyrosequencing

97 Sewage and sludge samples for pyrosequencing were taken from a full-scale MBR. The
98 reactor, as schematically shown in [Fig. S1 of the Supporting Information](#), was located in a general
99 merchandise store (31.3°N 121.4°E) of Shanghai, China and has been in operation for over 6 years.
100 The influent wastewater of the MBR includes 1) fresh food processing (FFP) wastewater, 2)
101 restaurant wastewater generated from restaurants serving Chinese, Japanese and Western style
102 food, 3) toilet flushing wastewater, 4) greywater from office region and washing basins, and 5) car
103 washing wastewater. The raw wastewater passed through screens, a dissolved air flotation tank
104 and an aerobic MBR tank. The treated wastewater was temporarily stored in an effluent tank and
105 finally reused for toilet flushing, lawn watering and car washing. The MBR tank had an effective
106 volume of 60 m³. 600 poly(vinylidene fluoride) flat-sheet membrane modules (Zizheng Environm
107 Technol Co. Ltd., Shanghai, China) with a mean pore size of 0.20 μm were installed in the tank.
108 Details about MBR setup, characteristics of the influent and treated wastewater are summarized in
109 the Supporting Information ([Section I, Fig. S1 and Table S1 in the Supporting Information](#)).

110 Influent wastewater, activated sludge and treated wastewater samples, termed as A1, A2 and
111 A3 samples, were taken from the inlet pipe, aerobic tank and outlet pipe of the MBR, respectively

112 (see Fig. S1). After DNA extraction and PCR amplification (see Section II of the Supporting
 113 Information), amplicons from A1, A2 and A3 were mixed at the equal concentration, and the
 114 mixture was used for pyrosequencing on a Roche 454 FLX Titanium platform at Majorbio
 115 Bio-Pharm Technology Co., Ltd (Shanghai, China).

116 2.2. Read quality control and subset construction

117 After pyrosequencing, 258,438 raw reads (0.1 G) were obtained according to the unique
 118 match to the barcodes (Table 1). The results were deposited into the NCBI short reads archive
 119 database (Accession Number: SRA169387). To improve the validity of subsequent data
 120 processing, Qiime (version 1.17 <http://qiime.org/>) was applied to (1) check the completeness of 3'
 121 end of primers and adaptors; (2) remove reads containing ambiguous base ('N') or homologous run
 122 that was longer than 10-nucleotide; (3) enable sliding window test of quality scores (-w 50 and -s
 123 20); and (4) remove reads shorter than 200 bps¹⁸. Barcodes and primers were also stripped from
 124 resulting sequences, and finally pyrosequencing produced 24,962 (A1), 113,131 (A2) and 54,525
 125 (A3) high-quality V1-V3 tags of the 16S rRNA gene with an average length of 462 bp (Table 1).

126 Table 1 Statistical summary for pyrosequencing and microbial diversity analysis

Sample ID	Raw reads	High-quality reads	Assigned reads ^a	OTU	Chao	Shannon	Fo
A1 ^b	34949	24962	19411	897	1323	3.70	0.92
B1	n.a. ^c	10000	8162	864	1886	3.88	0.89
C1	n.a.	100	89	40	164	2.90	0.64
A2	158938	113131	63243	1712	2026	4.85	0.90
B2	n.a.	10000	6163	1132	2664	5.04	0.81
C2	n.a.	100	63	46	187	3.65	0.42
A3	64551	54525	36644	1362	1670	5.20	0.89
B3	n.a.	10000	7132	1063	2092	5.34	0.83
C3	n.a.	100	75	59	190	3.95	0.37

127 a. Assigned reads are the reads that match the OTU in each sample. Some high-quality reads may not match any
 128 OTU for these reasons: (1) the read is chimeric, and (2) the read that has a singleton sequence is discarded.

129 b. A1, A2 and A3 represent pyrosequencing results of influent wastewater, activated sludge and treated
 130 wastewater samples; B1, B2 and B3 represent the subsets with 10000 reads extracted from A1, A2 and A3; C1, C2
 131 and C3 represent the subsets with 100 reads extracted from B1, B2 and B3.

132 c. n.a. indicates the value is not available.

133 For a comprehensive understanding of the impacts of sequencing depth, subsets with the
 134 capacity of 10,000 and 100 were generated from the high-quality maternal sets of A1~A3 by a
 135 semi-random extraction method. Initially, the sub.samlpe command of MOTHR program
 136 (<http://www.mothur.org/wiki/Sub.sample>) was used for A1~A3 to create 30 subsets comprised of

137 10,000 sequences, i.e., $B1_i = \{x_j | x_j \in A1, j = 1 \sim 1000\}$ ($i = 1 \sim 10$), $B2_i = \{x_j | x_j \in A2, j = 1 \sim 1000\}$
138 ($i = 1 \sim 10$) and $B3_i = \{x_j | x_j \in A3, j = 1 \sim 1000\}$ ($i = 1 \sim 10$). Principal coordinates analysis (PCoA)
139 with the Bray-Curtis index (R package, <http://www.r-project.org/>) was then performed to evaluate
140 the relationship between $A1 \sim A3$ and $B1_i \sim B3_i$ ($i = 1 \sim 10$), and the subsets with the highest
141 homology were retained and specified as $B1 \sim B3$ (Fig. S2 in the Supporting Information).
142 Afterwards, a similar procedure was applied to create 30 subsets containing 100 sequences from
143 $B1 \sim B3$, and the subsets with the highest homology with $B1 \sim B3$ were specified as $C1 \sim C3$. Despite
144 the debate that semi-random extraction is reliable enough compared to independent sequencing,
145 this method is similar to pyrosequencing run in reverse; and in practice the final gene libraries
146 (e.g., $A1 \sim A3$) can be obtained based on the deficient datasets (e.g., $B1 \sim B3$) by further sequencing
147 of the amplicons. Nevertheless, evaluation of pathogenic bacteria diversity was mainly based on
148 the original pyrosequencing results.

149 2.3. Phylogenetic classification and biodiversity analysis

150 Cluster of the high-quality reads into operational taxonomic units (OTUs) was performed
151 using UPARSE pipeline (vsesion 7.1, <http://drive5.com/uparse/>)¹⁹. Briefly, abundance-sorted
152 reads of the nine datasets ($A1 \sim A3$, $B1 \sim B3$ and $C1 \sim C3$) were clustered by setting a minimum
153 identity of 97%, and the uchime_ref command was used to filter out chimeras. The abundances of
154 OTUs in each dataset were obtained by searching the reads as a query set against the OTU
155 representative sequences. For the cluster files, alpha-diversity and rarefaction curves were
156 generated in MOTHUR for each sample (version v.1.30.1, <http://www.mothur.org>). Functional
157 organization indices (F_o) were calculated according to the standard method reported by Marzorati
158 et al.²⁰. Representative sequences from each OTU were assigned down to the phylum and genus
159 level using the RDP Classifier with a set confidence threshold of 80%
160 (<https://rdp.cme.msu.edu/classifier/classifier.jsp>, 16S rRNA training set 10)¹⁷.

161 Venn diagrams with shared and unique OTUs were utilized to depict the similarity and
162 difference between microbial communities. A pairwise statistical comparison of taxonomy at
163 phylum level between maternal sets and subsets was carried out using STAMP²¹. Biological
164 relevance between samples at genus level was evaluated using linear regression of SigmaPlot
165 software (version 12.5, Systat Software, Inc., U.S.). Furthermore, LDA Effect Size (LEfSe)
166 algorithm was introduced herein to identify taxa that characterize the differences among the three

167 environmental samples²². A1~A3, B1~B3 and C1~C3 were grouped according to the source (e.g.,
168 influent wastewater, activated sludge or treated wastewater sample), and each sample was firstly
169 normalized to the sum of the values of 0.05 M. The parameters for data processing were set as
170 follows: 'alpha value for the factorial Kruskal-Wallis test among classes' = 0.05, 'threshold on the
171 logarithmic LDA score for discriminative features' = 2.7, and 'set the strategy for multi-class
172 analysis' = all-against-all.

173 Alignment of microbial communities to pathogenic genera was firstly evaluated using the
174 taxonomic results of RDP Classifier. The lists of known pathogenic genera summarized by [Ye and](#)
175 [Zhang](#) and [Bidy et al.](#) were used as reference^{7, 23}. Representative sequences from OTUs that
176 were assigned into *Arcobacter*, *Clostridium*, *Legionella*, and *Mycobacterium* were further
177 separated for phylogenetic analysis at species level. 16S rRNA gene of known pathogens and
178 non-pathogens from the four genera were achieved from NCBI Genbank ([Table S2 of the](#)
179 [Supporting Information](#)), and merged with the corresponding sequences of this study into a fasta
180 file. ClustalW was used for aligning and bootstrapping of the phylogenetic tree, which was then
181 viewed, edited and published with MEGA 6²⁴. Default settings were used. Furthermore, bacteria
182 assigned to the families of Enterobacteriaceae and Enterococcaceae were regarded as the
183 representative fecal indicators in this study.

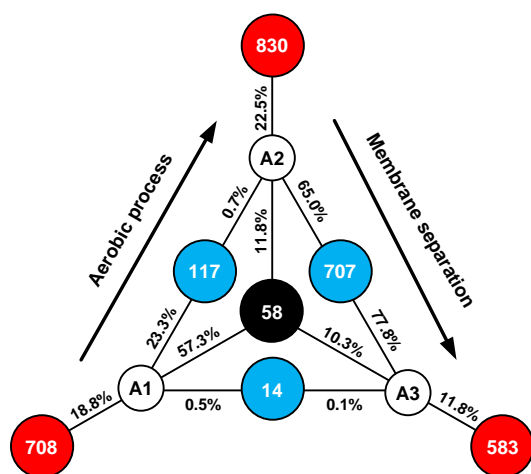
184 2.4. Quantification of bacterial biomass using flow cytometer (FCM)

185 Bacteria biomass in wastewater and sludge samples was quantified using flow cytometer.
186 Initially, influent wastewater, activated sludge and treated wastewater (A1, A2 and A3) were
187 diluted 1:20 (v/v), 1:500 (v/v) and 1:1 (v/v) using 0.22- μ m filtered phosphate-buffered-saline
188 solution (0.84 %, pH = 7) to achieve optimal concentrations of bacteria for FCM analysis. Then
189 the diluted mixtures (A1, A2 and A3) were subjected to ultrasonication treatment at power
190 densities of 25, 80 and 0 kJ/L, respectively. After filtrated with 10- μ m filters, samples were
191 stained with SYBR Green I at a ratio of 100:1, incubated for 15 min in the dark at room
192 temperature and finally processed to the flow cytometer (BD AccuriTM C6, U.S.). Each sample
193 was test in triplicate and total bacterial counts of A1, A2 and A3 were $(2.31 \pm 0.24) \times 10^8$, $(7.06 \pm$
194 $0.30) \times 10^9$ and $(3.35 \pm 0.82) \times 10^4$ counts/mL, respectively.

195 3. Results

196 3.1. Diversity and similarity analysis of microbial communities

197 By performing the alignment at α of 0.03 using UPARSE pipeline, 897, 1712 and 1362 OTUs
 198 were obtained from A1, A2 and A3 (Table 1). At a degraded and uniform library size of 10,000,
 199 Chao1 richness estimators of the three samples were 1886, 2664 and 2092, and Shannon diversity
 200 indices were 3.88, 5.04 and 5.34, respectively. Alpha-diversity analysis suggested that bacterial
 201 community from the influent wastewater sample had the lowest microbial richness and diversity.
 202 Moreover, F_o of the three samples were 0.89~0.92. It could be inferred that all the microbial
 203 communities were highly functionally organized^{14,20}. Pairwise comparison using Venn analysis
 204 showed that the similarity of A2-A3 was the highest, followed by that of A1-A2 and that of
 205 A1-A3 (Fig. 1). Notably, A2 and A3 had 707 shared OTUs that contained 65.0% and 77.8% of the
 206 reads, respectively. In contrast, only very few reads (0.5% and 0.1%) were classified into the
 207 OTUs that were shared by A1 and A3, which indicated that MBR treatment introduced a profound
 208 influence on the structure of microbial community in wastewater.



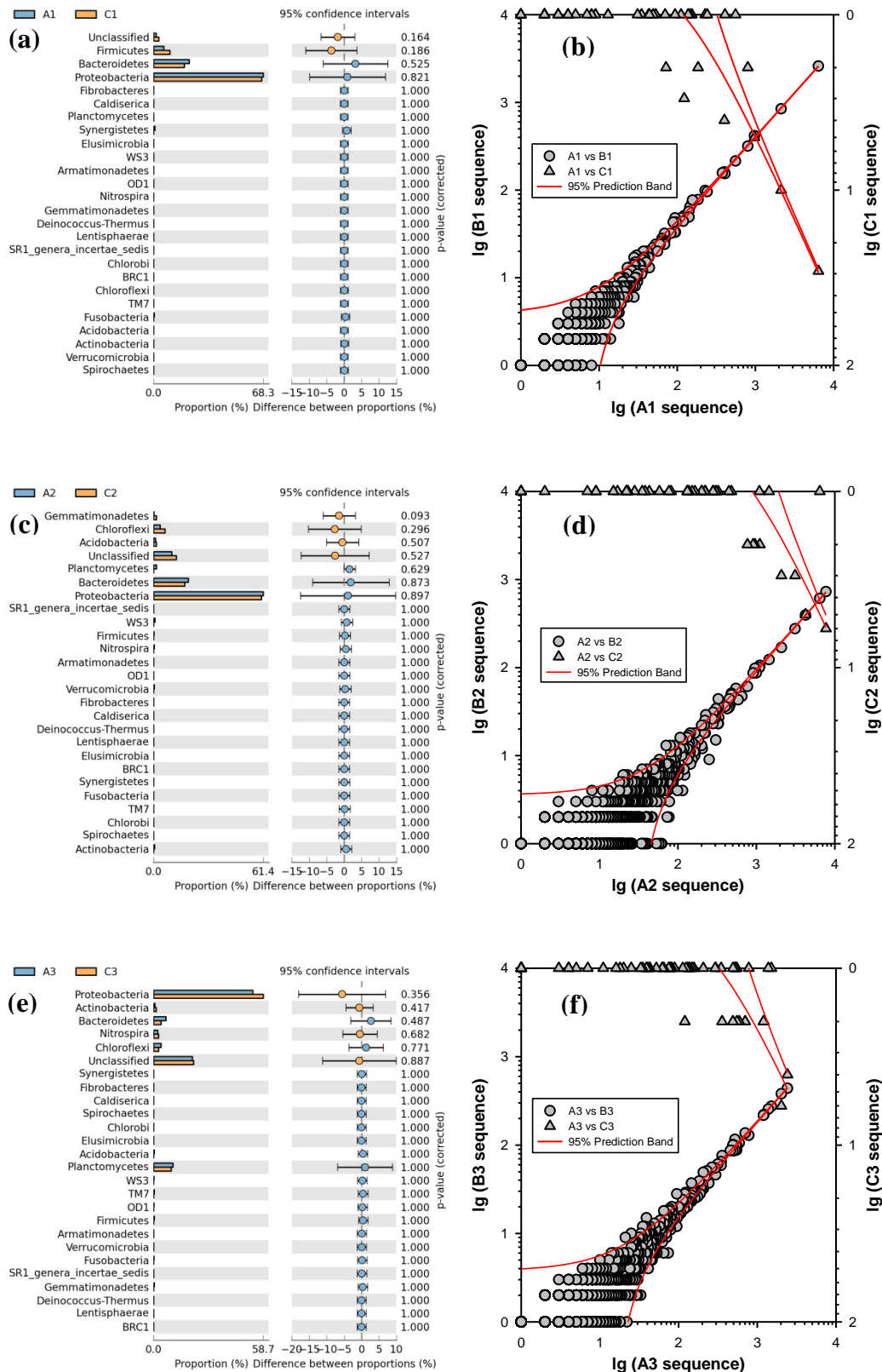
209
 210 Fig. 1 Similarity analysis of the microbial communities (A1, A2 and A3) based on the clustering
 211 results at 3% distance cutoff. The numbers in the black circles represent the number of OTUs that
 212 is present in the core OTUs shared by the three samples. The numbers in the blue circles represent
 213 the OTUs shared by two samples. The numbers in the red circles represent the unique OTUs
 214 observed in only one sample. Percentages listed beside the branches indicate the percentages of
 215 reads of each sample assigned into the nearby OTUs groups.

216 3.2. Impacts of sequencing depth

217 Generally, microbial communities of environmental samples are highly diverse, and in this
 218 study rarefaction curves showed that new bacterial phylotypes continued to emerge even after
 219 60,000 reads sampled (Fig. S3 in the Supporting Information). Addressing an appropriate

220 sequencing depth is crucial for high-throughput pyrosequencing to detect pathogenic bacteria at
221 low abundance; an enlarged depth significantly increases the sequencing and processing cost,
222 while a small library size could only provide insufficient resolution. Since 10,000 and 100 library
223 sizes are always considered in pyrosequencing and conventional molecular biology studies,
224 sequence subsets with corresponding capacities were generated from the maternal gene libraries
225 by the semi-random extraction method.

226 [Table 1](#) indicated that insufficient resolution reduced the accuracy of alpha-diversity analysis.
227 For instance, at the sequencing depth of 100, only 40, 46 and 59 OTUs were predicted for the
228 whole microbial communities in influent wastewater, activated sludge and treated wastewater,
229 respectively. Furthermore, we compared the taxonomic results of the nine datasets at phylum and
230 genus levels ([Fig. 2](#)). In total, 25 phyla were classified at the threshold of 80%. Proteobacteria was
231 the most dominant phylum, accounting for 53.1~68.3% of total communities, respectively.
232 Pairwise comparison using STAMP shows that there is no significant dissimilarity of taxonomic
233 results between A_j and C_j ($j = 1, 2, 3$) at phylum level ([Fig. 2a](#), [Fig. 2c](#) and [Fig. 2e](#)). However, the
234 reliability was significantly declined with the taxonomy down to genus level. Most taxa of A1~A3
235 could not be predicted by the taxonomic results of C1~C3 at a 95% predication band. In contrast,
236 linear regression showed that except for a few categories, B1~B3 supplied a credible
237 characterization of the microbial communities of A1~A3 at genus level ([Fig. 2b](#), [Fig. 2d](#) and [Fig.](#)
238 [2f](#)). The results suggested that compared to low-throughput sequencing methods pyrosequencing
239 could provide a more valid estimation of the population structure of diverse communities
240 especially at terminal taxonomic levels (e.g., genus level). Increasing the library size from 10,000
241 to 100,000, however, did not improve the taxonomic results as expected, probably due to the
242 abundance of singletons at high sequencing depth, which were always discarded after denoising
243 ([Table 1](#)).



244

245

246

247

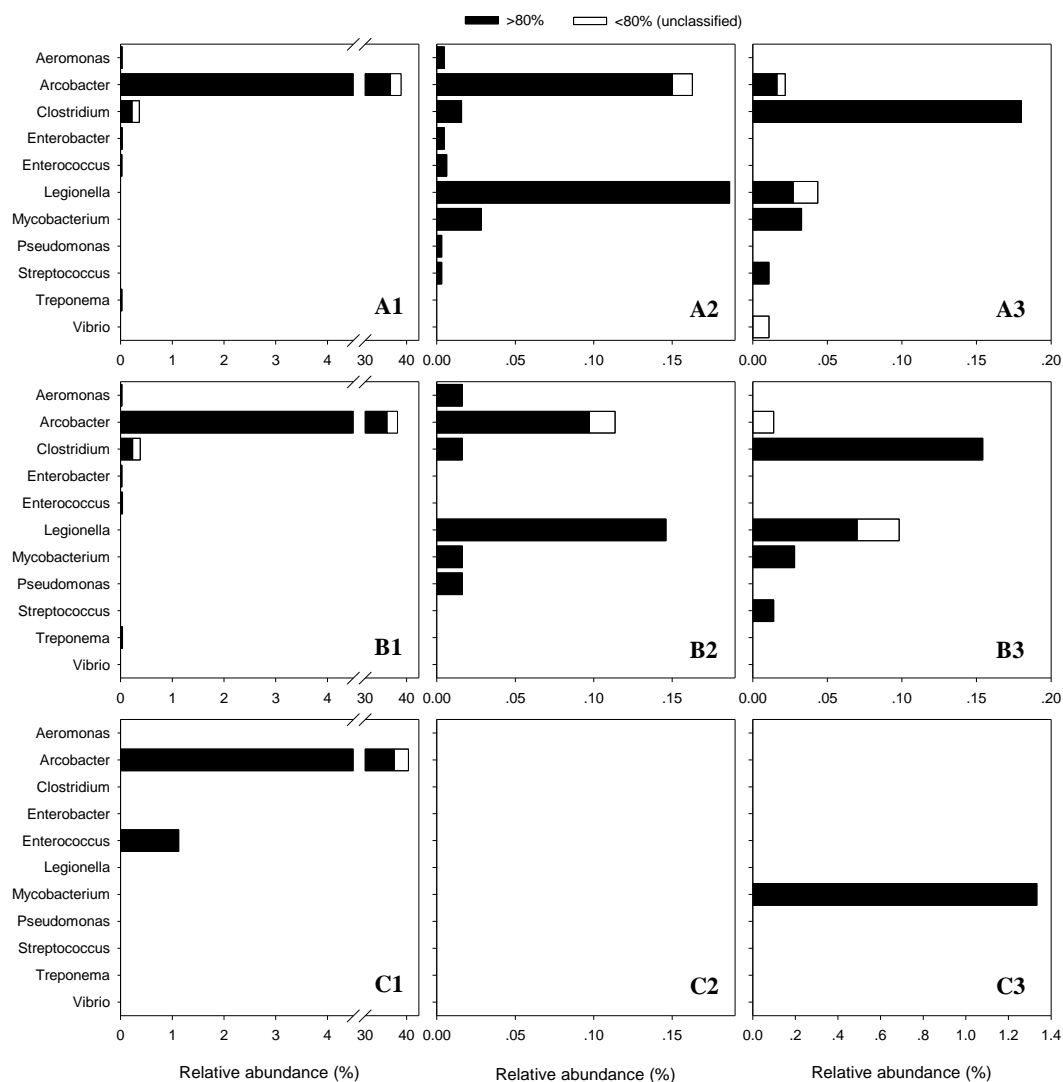
248

249 Fig. 2 Pairwise comparison of biological relevance of (a) A1 and C1 at phylum level, (b) A1, B1
 250 and C1 at genus level, (c) A2 and C2 at phylum level, (d) A2, B2 and C2 at genus level, (e) A3
 251 and C3 at phylum and (f) A3, B3 and C3 at genus level. Taxonomic results based on OTU
 252 clustering at a 3% distance were compared using STAMP at phylum level. A corrected P -value

253 lower than 0.05 is significant. Correlations of assignment results in each of the three samples were
254 carried out at genus level. The horizontal and vertical axes in each subfigure (b, d and f) indicate
255 the numbers of the corresponding genus sequences. The red lines represent the 95% prediction
256 bands of linear regression.

257 3.3. Detection and characterization of the potential pathogenic bacteria

258 During MBR treatment, the structure of microbial community in wastewater changed in
259 response to the environmental selective pressures, and taxa were differently enriched in different
260 samples (Fig. S4 of the Supporting Information). Fig. 3, according to the alignment to the lists of
261 known pathogenic genera ^{7,23}, shows the eleven genera of potential pathogenic bacteria found in
262 the three samples. It could be noticed that only *Arcobacter* and *Clostridium* were ubiquitous in all
263 the samples. In A1, *Arcobacter* was the most abundant genera, accounting for about 40% of the
264 population. Except for *Clostridium*, the other potential pathogens, including *Aeromonas*,
265 *Enterobacter*, *Enterococcus* and *Treponema*, were present at very low abundances
266 (0.026~0.031%). The number of sequences assigned into potential pathogenic bacteria was
267 significantly decreased in the activated sludge sample (Table S3 of the Supporting Information).
268 For example, *Arcobacter* were underrepresented in A2, with nearly two orders of magnitude
269 difference in abundance compared to those found in A1. A similar decay was found for
270 *Aeromonas*, *Enterobacter*, *Enterococcus* and *Treponema* at even lower abundances. Instead, a
271 well-known source of infection, *Legionella* ²⁵, became abundant among the potential pathogens.
272 Furthermore, membrane retention induced a selective pressure on bacterial pathogen diversity.
273 Several widely-reported pathogenic bacteria, including *Aeromonas*, *Enterobacter*, *Enterococcus*,
274 and *Pseudomonas* ^{7, 8, 23, 26}, were not found in A3, and the sequences assigned to potential
275 pathogenic genera only accounted for 0.3% of the total population (Table S3). However, a
276 gram-positive genus, *Clostridium*, was found to be the most abundant among the potential
277 pathogenic phylotypes. Since no significant difference was noted by LEfSe analysis at this taxon
278 (Fig. S4), it was possible to infer that *Clostridium* were more resistant to the treatment of MBR.

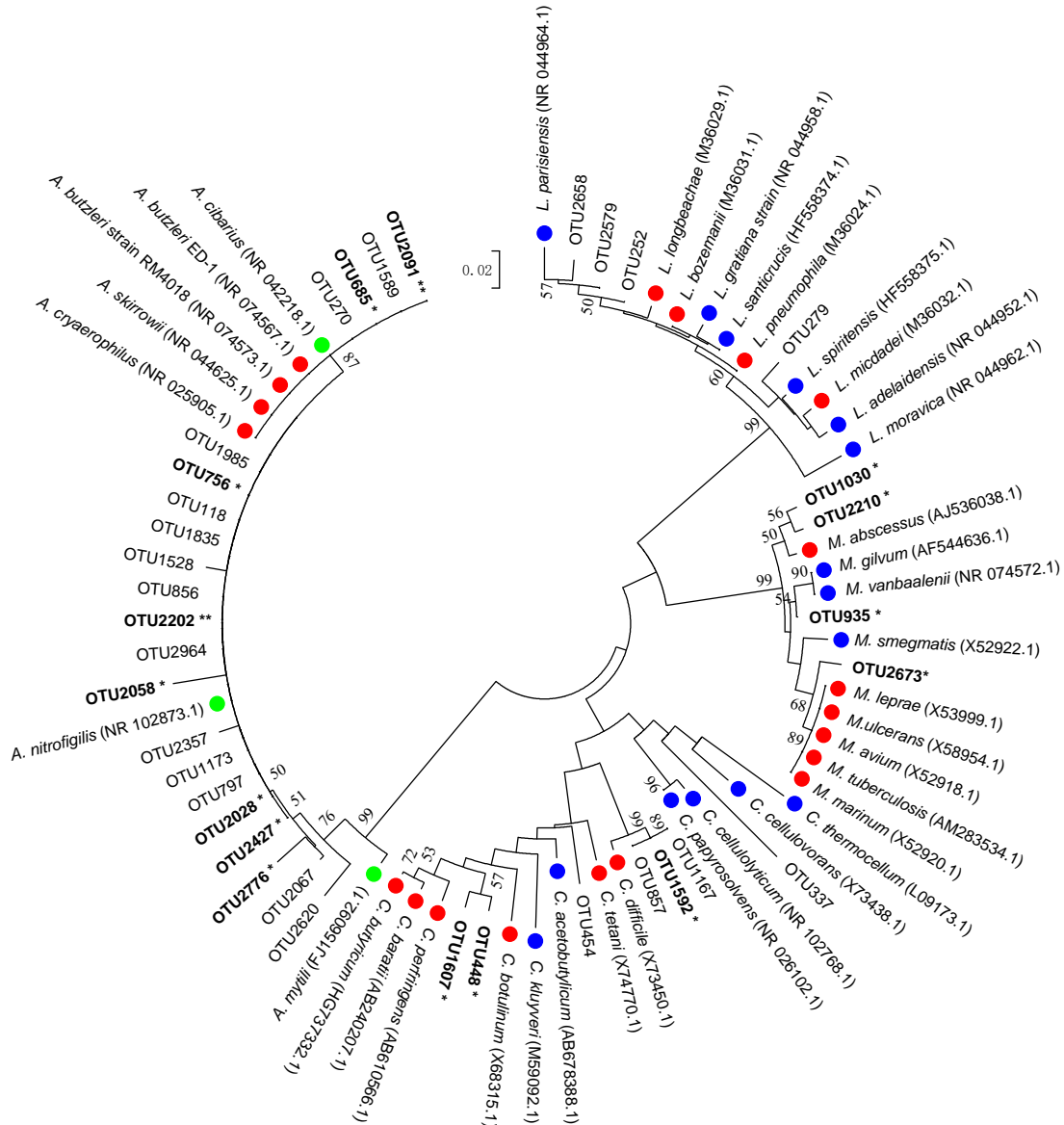


279

280 Fig. 3 Relative abundances of potential pathogenic genera in influent wastewater (A1), activated
 281 sludge (A2), treated wastewater (A3), subsets with 10000 reads (B1, B2 and B3) and subsets with
 282 100 reads (C1, C2 and C3). Relative abundance is defined as the percentage of a pathogenic genus
 283 in total population. The different colors represent the percentages of sequences in the
 284 corresponding confidence ranges.

285 In this study, the depth of pyrosequencing also had a significant influence on detection and
 286 characterization of potential pathogenic bacteria. For B1~B3, the dominant pathogenic genera
 287 could be well identified, while those with low abundances (e.g., *Enterobacter*, *Enterococcus* and
 288 *Vibrio*) were neglected (Fig. 3). It is worthy noting that sequencing failed to reveal the majority of
 289 potential pathogenic bacteria in environmental samples when the library size was declined to 100.
 290 Specifically, no pathogenic bacteria was detected in C2 and *Mycobacterium* were obviously
 291 overestimated in C3 (Fig. 3).

292 Overall, the corresponding sequences got a good alignment with the potential pathogenic
293 genera by using the RDP Classifier, and most had a bootstrap confidence over 80% (Fig. 3 and
294 Table S3). Since species-specific comparison with known pathogenic bacteria could give a more
295 accurate estimation of the potential pathogens in the samples, representative sequences of
296 concerned OTUs assigned into *Arcobacter*, *Clostridium*, *Legionella* and *Mycobacterium* genera
297 (Fig. 3) were retrieved from the datasets. Phylogenetic analysis was then conducted by building a
298 library with representative 16S rRNA gene sequences of pathogenic and non-pathogenic bacterial
299 species. As shown in Fig. 4, 15 OTUs of the total (37 OTUs) had an identity over 95% with
300 neighbor pathogens, including 8 OTUs assigned into *Arcobacter*, 3 OTUs into *Clostridium* and 4
301 OTUs into *Mycobacterium* (Table S4 of the Supporting Information). In *Arcobacter* genus, the
302 most abundant taxon OTU2964, accounting for 32.9% of the population in A1, showed a low
303 alignment with known pathogens. Without additional information, it could not be concluded
304 whether these sequences referred to nonpathogenic strains, because variants might be also
305 associated with disease but not yet identified²³. Of particulate importance is that only two strains
306 (OTU2202 and OTU2091) were recognized as potential pathogenic species by species-specific
307 comparison. It could be deduced that phylogenetic analysis at the genus level might lead to an
308 overestimation of the pathogenic bacteria in environmental samples.



309

310 Fig. 4 Phylogenetic tree of concerned OTUs from A1~A3 that were assigned into potential
 311 pathogenic genera (*Arcobacter*, *Clostridium*, *Legionella* and *Mycobacterium*). Representative
 312 sequences from these OTUs were reahieved for alignment and phylogenetic analysis. The number of
 313 OTU (e.g., OTU448) only indicates the logical order in OTU clustering. OTUs with satisfied
 314 identities to neighbor pathogens are bolded and marked with * (95~99%) or ** (>99%). Bootstrap
 315 values are calculated by 1000 repetitions, and values >50% are given. ● indicate the
 316 non-pathogenic species, ● the pathogenic species and ● the vague species.

317 Since fecal indicator bacteria are still widely used to predict the presence of bacterial, viral
 318 and protozoan pathogens⁹, the abundances of typical indicators were further evaluated based on
 319 the taxonomic results in the present work. Because there is no full taxonomic definition of fecal

320 indicator yet, bacteria assigned to the families of Enterobacteriaceae and Enterococcaceae were
 321 regarded as the representative fecal indicators herein. As shown in Table 2, 41 sequences from A1
 322 were classified, which contributed 0.21% of the dataset. In A2, only 10 sequences got a valid
 323 match by RDP Classifier, including 6 sequences assigned into Enterobacteriaceae and 4 sequences
 324 into Enterococcaceae. Probably due to the sound separation of 0.20- μ m poly(vinylidene fluoride)
 325 membranes, no Enterobacteriaceae or Enterococcaceae was detected in the treated wastewater (A3)
 326 in this study.

327 Table 2 Summary of sequences assigned to Enterobacteriaceae and Enterococcaceae ^a

	A1		A2		A3	
	Number of sequences	<i>r</i> , % ^b	Number of sequences	<i>r</i> , %	Number of sequences	<i>r</i> , %
Enterobacteriaceae	36	0.185	6	0.009	0	0
Enterococcaceae	5	0.026	4	0.006	0	0
Total	41	0.211	10	0.015	0	0

328 a. based on the taxonomic results of RDP Classifier

329 b. *r* indicates the relative abundance of sequences.

330 4. Discussion

331 In this study, a group of predominant potential pathogens, *Arcobacter*, were differently
 332 abundant in the influent wastewater compared to other samples (Fig. 3 and Fig. S4). The genus
 333 *Arcobacter* belonging to the RNA Superfamily VI of Proteobacteria was proposed in 1991, and
 334 the International Commission on Microbiological Specification for Foods has considered
 335 *Arcobacter* to be one of the most frequently notified food-borne infectious agents ²⁷. Full
 336 understanding of its occurrence and fate during wastewater reclamation is, thereby, very important,
 337 especially for a rapid and accurate diagnosis of the infection source of outbreaks (e.g., acute
 338 enteric disease) in public places. In the present work, RDP Classifier indicated that 21 OTUs from
 339 A1 were classified into *Arcobacter* genus, which accounted for 36.5% of the total population. The
 340 microbial composition might be a typical pathogenic characteristic of restaurant wastewater,
 341 because raw and undercooked meat and poultry products have been recognized as the sources of
 342 *Arcobacter* ²⁸. Overall, the results showed that the hybrid MBR system presented a good removal
 343 of *Arcobacter*; ~87% of the influent *Arcobacter* were eliminated in the activated sludge and <10
 344 counts/mL were detected in the treated wastewater (Section III of the Supporting Information).
 345 Furthermore, phylogenetic analysis suggested that the genus-specific comparison could result in an

346 overestimation of pathogenic bacteria, since more than 90% of *Arcobacter* had highest homology
347 with a free-living nitrogen-fixing bacterium, *A. nitrofigilis*. Notably, 2 strains (OTU2091 and
348 OTU2202) had an identity over 99% with known pathogenic *Arcobacter* species (e.g., *A. butzleri*),
349 which were estimated to be at concentrations of $(1.40 \pm 0.15) \times 10^6$, $(6.81 \pm 0.29) \times 10^6$ and ~ 1
350 counts/mL in influent wastewater, activated sludge and treated wastewater, respectively (Section
351 III of the Supporting Information).

352 In the view of pathology, it is of great concern to focus on the pathogenic bacteria emerging
353 in the aerobic tank of MBR because not only these tolerant microorganisms got a competitive
354 advantage with the biomarkers involved in contaminant degradation (e.g., *Zoogloea* and
355 *Dechloromonas* as shown in Fig. S4) but also aerosols containing pathogens could be generated
356 from the aeration tank and further transported and dispersed by wind. A genus of gram-negative
357 coccobacilli, *Legionella*, was well recognized in the activated sludge sample (A2). It has been
358 reported that *Legionella* prefer to inhabit in man-made aquatic environments where the water
359 temperature is higher than ambient temperature, and that the growth of *Legionella* spp. can be
360 aided by co-existing micro-organisms (e.g., protozoa)^{25, 29}. Although this bacterial genus was
361 enriched in activated sludge, *Legionella* found herein seemed non-pathogenic; all the
362 representative sequences had a low alignment (87~92%) with the foremost pathogenic species *L.*
363 *pneumophila*, *L. longbeachae*, *L. micdadei* and *L. bozemanii*. Fig. 3 indicates that MBR removal
364 of *Legionella* from restaurant wastewater was mainly attributed to membrane retention that could
365 efficiently eliminate the hosts (e.g., amoebae) in the treated wastewater. Moreover, a recent study
366 on bacterial pathogen diversity in biosolids (digested sludge) using pyrosequencing has revealed
367 that most of the pathogenic sequences belonged to the genera of *Mycobacterium* and *Clostridium*
368²³. In the present work, our results showed that despite low relative abundances, all sequences
369 belonging to *Mycobacterium* genus had more than 95% similarity to a 'freak' pathogenic species,
370 *M. abscessus*. The gene order phylogeny of *M. abscessus* groups the organisms into rapid and
371 slow-growers^{30, 31}. *M. abscessus* is closer to the non-pathogens in terms of its growth
372 characteristics and is placed away from the pathogens (Fig. 4), which could lead to taxonomic bias
373 based on 16S rRNA gene pyrosequencing. As a result, virulence assays that target the functional
374 genes are further required to convince the relevant conclusions.

375 This study also reinvigorates the debate that the indicator bacteria are inefficient in
376 representing the potential pathogenic bacteria from non-fecal sources. The relative abundance of
377 fecal indicator bacteria did not show a good relationship with that of potential pathogenic species,
378 though providing a rough evaluation on the occurrence of potential pathogens in the restaurant
379 wastewater (Table 2 and Table S4). Notably, Enterobacteriaceae and Enterococcaceae were not
380 detected in A3 but 97 of total 36644 sequences were classified into the pathogenic genera,
381 including 25 sequences with identities of > 95% with *M. abscessus*, *C. difficile* and *C. botulinum*.
382 Furthermore, *Clostridium* genus was an important group in A3, which was resistant to the MBR
383 treatment (Fig. 3). In this study, the ambiguously defined taxon contained *Clostridium* cluster
384 sensu stricto (Clostridiaceae 1), *Clostridium* cluster IV (Ruminococcaceae) and *Clostridium*
385 cluster XI (Peptostreptococcaceae). Three strains from *Clostridium* cluster XI and *Clostridium*
386 cluster sensu stricto had a good phylogenetic alignment with *C. difficile* and *C. botulinum*,
387 respectively (Fig. 4 and Table S4). Since *Clostridia* (spores) are highly resistant to chlorination³²,
388 disinfection efficiency could be easily overestimated when referred to the elimination of intolerant
389 indicator bacteria (e.g., Enterobacteriaceae).

390 In this study, 454 pyrosequencing was introduced for a comprehensive understanding of
391 bacterial pathogens in the restaurant wastewater. Compared to conventional culture-based
392 methods and qPCR assays, this technology is high-throughput for mining potential pathogenic
393 bacteria in environmental samples, which avoids the misestimation of pathogens by using a certain
394 group of indicator bacteria. Molecular biology methods that offer ≈ 100 tags could only provide
395 rough information on the structure of microbial communities at phylum level (Fig. 2). By contrast,
396 the 10,000-sequence datasets were generally valid in forecasting individuals within microbial
397 communities, but overrepresentation and underrepresentation were still noted regarding the
398 highly-potential bacterial pathogens (Table S4). For accurately reaping the rare strains,
399 exponential growth of the library size (1~2 orders of magnitude) might be unwise since a large
400 number of singletons were generated at a sequencing depth of 30,000~150,000 (Table 1). DNA
401 fragment pretreatment (e.g., the use of multiple genus level PCR primers) should be thereby
402 considered in pathogenic studies. Furthermore, 16S rRNA gene pyrosequencing provided the
403 opportunity to discover the important strains that have not been cultured yet (e.g., OTU2964).
404 Short-gun metagenomic and metatranscriptomic sequencing could be used to predict their

405 functions. Phylogenetic analysis of concerned pathogenic and non-pathogenic bacterial species
406 herein gave a more accurate evaluation of the abundance and diversity of bacterial pathogens.
407 Virulence of relevant communities could be further analyzed using qPCR or microarrays that
408 target the functional genes. Overall, the present work showed that restaurant wastewater was
409 suitable for reclamation using MBR technology. Pathogenic bacteria were efficiently removed in
410 the hybrid systems, and membrane filtration process retained the communities that were resistant
411 to biological treatment. The tolerant bacterial pathogens in treated wastewater revealed by
412 pyrosequencing provide insights into the selection of specific tertiary treatment and proper
413 disinfection methods as well.

414 5. Conclusions

415 In the present work, high-throughput pyrosequencing was used to characterize the potential
416 pathogenic bacteria in a full-scale MBR treating restaurant wastewater. The results indicated that
417 the influent pathogenic community might be highly diverse and that 39.2% of the population was
418 assigned into the pathogenic genera. Overall, MBR treatment had a good removal of *Aeromonas*,
419 *Arcobacter*, *Enterobacter*, *Enterococcus* and *Treponema*, and in the treated wastewater the
420 bacterial count of *Arcobacter* was decreased to <10 counts/mL. The most dominant potential
421 pathogens in activated sludge and treated wastewater were affiliated to the genera of *Legionella*,
422 *Clostridium* and *Mycobacterium*. Nevertheless, species-specific comparison showed that only a
423 small portion (0.0~1.6%) of the corresponding sequences had identities of > 99% to the neighbor
424 pathogenic species, suggesting that phylogenetic analysis at the genus level might lead to an
425 overestimation of the potential pathogens. This study provided insights into assessing pathogenic
426 bacteria risk in wastewater purification and reclamation.

427

428 Acknowledgments

429 The work is financially supported by National Natural Science Foundation of China
430 (51422811), Shanghai Rising-Star Program (14QA1403800) and the Shanghai Science &
431 Technology Commission Program (13231202002).

432

433 References

434 1. Chinese catering information center. <http://www.canyin.com/a/zixunzhongxin>

- 435 /xingyexinwen/20120724/991.html. (accessed April 18, 2014).
- 436 2. X. Zhu, Z. Wang and Z. Wu, *Process Biochem.* , 2011, **46**, 1001-1009.
- 437 3. B. Yang, G. Chen and G. Chen, *Sep. Purif. Technol.* , 2012, **88**, 184-190.
- 438 4. J.-x. Kang, L. Lu, W. Zhan, B. Li, D.-S. Li, Y.-Z. Ren and D.-q. Liu, *J. Hazard. Mater.* ,
439 2011, **186**, 849-854.
- 440 5. S. Judd, *Trends Biotechnol.* , 2008, **26**, 109-116.
- 441 6. M. E. Verbyla, S. M. Oakley and J. R. Mihelcic, *Environ. Sci. Technol.*, 2013, **47**,
442 3598-3605.
- 443 7. L. Ye and T. Zhang, *Environ. Sci. Technol.*, 2011, **45**, 7173-7179.
- 444 8. W. Ahmed, H. Brandes, P. Gyawali, J. P. S. Sidhu and S. Toze, *Water Res.* , 2014, **53**,
445 361-369.
- 446 9. O. Savichtcheva and S. Okabe, *Water Res.* , 2006, **40**, 2463-2476.
- 447 10. W. Ahmed, S. Sawant, F. Huygens, A. Goonetilleke and T. Gardner, *Water Res.* , 2009, **43**,
448 4918-4928.
- 449 11. K. Zhang and K. Farahbakhsh, *Water Res.* , 2007, **41**, 2816-2824.
- 450 12. P. Bourhy, S. Bremont, F. Zinini, C. Giry and M. Picardeau, *J. Clin. Microbiol.* , 2011, **49**,
451 2154-2160.
- 452 13. D.-Y. Lee, H. Lauder, H. Cruwys, P. Falletta and L. A. Beaudette, *Sci. Total Environ.* ,
453 2008, **398**, 203-211.
- 454 14. J. Ma, Z. Wang, Y. Yang, X. Mei and Z. Wu, *Water Res.* , 2013, **47**, 859-869.
- 455 15. T. Zhang, M. F. Shao and L. Ye, *Isme Journal*, 2012, **6**, 1137-1147.
- 456 16. A. S. Laufer, J. P. Metlay, J. F. Gent, K. P. Fennie, Y. Kong and M. M. Pettigrew, *Mbio*,
457 2011, **2**.
- 458 17. Q. Wang, G. M. Garrity, J. M. Tiedje and J. R. Cole, *Appl. Environ. Microbiol.*, 2007, **73**,
459 5261-5267.
- 460 18. J. G. Caporaso, J. Kuczynski, J. Stombaugh, K. Bittinger, F. D. Bushman, E. K. Costello,
461 N. Fierer, A. G. Pena, J. K. Goodrich, J. I. Gordon, G. A. Huttley, S. T. Kelley, D. Knights,
462 J. E. Koenig, R. E. Ley, C. A. Lozupone, D. McDonald, B. D. Muegge, M. Pirrung, J.
463 Reeder, J. R. Sevinsky, P. J. Tumbaugh, W. A. Walters, J. Widmann, T. Yatsunenko, J.
464 Zaneveld and R. Knight, *Nat. Methods* 2010, **7**, 335-336.

- 465 19. R. C. Edgar, *Nat. Methods* 2013, **10**, 996-+.
- 466 20. M. Marzorati, L. Wittebolle, N. Boon, D. Daffonchio and W. Verstraete, *Environ.*
467 *Microbiol.* , 2008, **10**, 1571-1581.
- 468 21. D. H. Parks and R. G. Beiko, *Bioinformatics*, 2010, **26**, 715-721.
- 469 22. N. Segata, J. Izard, L. Waldron, D. Gevers, L. Miropolsky, W. S. Garrett and C.
470 Huttenhower, *Genome Biol.*, 2011, **12**.
- 471 23. K. Bibby, E. Viau and J. Peccia, *Water Res.* , 2010, **44**, 4252-4260.
- 472 24. K. Tamura, G. Stecher, D. Peterson, A. Filipinski and S. Kumar, *Mol. Biol. Evol.*, 2013, **30**,
473 2725-2729.
- 474 25. J. M. Blatny, B. A. P. Reif, G. Skogan, O. Andreassen, E. A. Hoiby, E. Ask, V. Waagen, D.
475 Aanonsen, I. S. Aaberge and D. A. Caugant, *Environ. Sci. Technol.*, 2008, **42**, 7360-7367.
- 476 26. H. de Man, M. Bouwknecht, E. van Heijnsbergen, E. J. T. M. Leenen, F. van Knapen and A.
477 M. de Roda Husman, *Water Res.* , 2014, **54**, 254-261.
- 478 27. I. Gonzalez, T. Garcia, S. Fernandez and R. Martin, *Food Anal. Methods* 2012, **5**,
479 956-968.
- 480 28. M. Mor-Mur and J. Yuste, *Food Bioprocess Technol.* , 2010, **3**, 24-35.
- 481 29. H. Y. Lau and N. J. Ashbolt, *J. Appl. Microbiol.* , 2009, **107**, 368-378.
- 482 30. H. Medjahed, J.-L. Gaillard and J.-M. Reyrat, *Trends Microbiol.*, 2010, **18**, 117-123.
- 483 31. A. N. Prasanna and S. Mehra, *Plos One*, 2013, **8**.
- 484 32. T. Karpova, P. Pekonen, R. Gramstad, U. Ojstedt, S. Laborda, H. Heinonen-Tanski, A.
485 Chavez and B. Jimenez, *Water Sci. Technol.* , 2013, **68**, 2090-2096.

486

487

488 **Figure captions**

489 **Fig. 1** Similarity analysis of the microbial communities (A1, A2 and A3) based on the clustering
490 results at 3% distance cutoff. The numbers in the black circles represent the number of OTUs that
491 is present in the core OTUs shared by the three samples. The numbers in the blue circles represent
492 the OTUs shared by two samples. The numbers in the red circles represent the unique OTUs
493 observed in only one sample. Percentages listed beside the branches indicate the percentages of
494 reads of each sample assigned into the nearby OTUs groups.

495 **Fig. 2** Pairwise comparison of biological relevance of (a) A1 and C1 at phylum level, (b) A1, B1
496 and C1 at genus level, (c) A2 and C2 at phylum level, (d) A2, B2 and C2 at genus level, (e) A3
497 and C3 at phylum and (f) A3, B3 and C3 at genus level. Taxonomic results based on OTU
498 clustering at a 3% distance were compared using STAMP at phylum level. A corrected *P*-value
499 lower than 0.05 is significant. Correlations of assignment results in each of the three samples were
500 carried out at genus level. The horizontal and vertical axes in each subfigure (b, d and f) indicate
501 the numbers of the corresponding genus sequences. The red lines represent the 95% prediction
502 bands of linear regression.

503 **Fig. 3** Relative abundances of potential pathogenic genera in influent wastewater (A1), activated
504 sludge (A2), treated wastewater (A3), subsets with 10000 reads (B1, B2 and B3) and subsets with
505 100 reads (C1, C2 and C3). Relative abundance is defined as the percentage of a pathogenic genus
506 in total population. The different colors represent the percentages of sequences in the
507 corresponding confidence ranges.

508 **Fig. 4** Phylogenic tree of concerned OTUs from A1~A3 that were assigned into potential
509 pathogenic genera (*Arcobacter*, *Clostridium*, *Legionella* and *Mycobacterium*). Representative
510 sequences from these OTUs were reahieved for alignment and phylogenic analysis. The number of
511 OTU (e.g., OTU448) only indicates the logical order in OTU clustering. OTUs with satisfied
512 identities to neighbor pathogens are bolded and marked with * (95~99%) or ** (>99%). Bootstrap
513 values are calculated by 1000 repetitions, and values >50% are given. ● indicate the
514 non-pathogenic species, ● the pathogenic species and ● the vague species.

515 **Table Captions**

516 **Table 1** Statistical summary for pyrosequencing and microbial diversity analysis

517 **Table 2** Summary of sequences assigned to Enterobacteriaceae and Enterococcaceae

518