# Machine Leaning-based Determination of Sampling Depth for Complex Environmental Systems: Case Study with Single-Cell Raman Spectroscopy Data in EBPR Systems

Guangyu Li[1,‡], Chieh Wu[2,‡], Dongqi Wang[3,1], Varun Srinivasan[1], David R. Kaeli[2], Jennifer G. Dy[2], April Z. Gu[4,*]

[1] Department of Civil and Environmental Engineering, Northeastern University, Boston, MA

[2] Department of Electrical and Computer Engineering, Northeastern University, Boston, MA

[3] Xi'an University of Technology, Xi'an, Shaanxi, PRC

[4] School of Civil and Environmental Engineering, Cornell University, Ithaca, NY
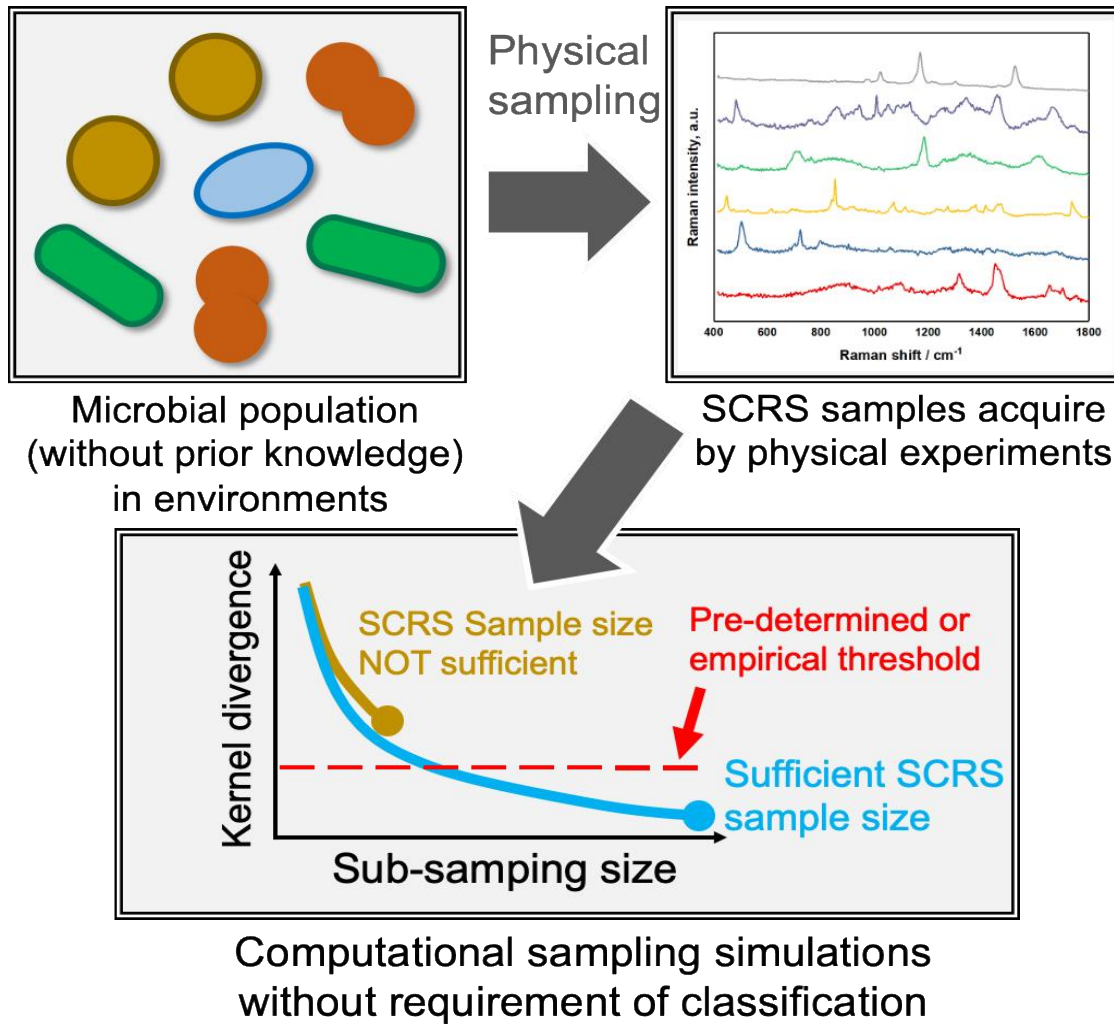
[‡] These authors contributed equally to this work

[*] Corresponding author: aprilgu@cornell.edu

**ABSTRACT:** Rapid progress in various advanced analytical methods such as single-cell technologies enable unprecedented and deeper understanding of microbial ecology beyond the resolution of conventional approaches. A major application challenge exists in the determination of sufficient sample size without sufficient prior knowledge of the community complexity and, the need to balance between statistical power and limited time or resources. This hinders the desired standardization and wider application of these technologies. Here, we proposed, tested and validated a computational sampling size assessment protocol taking advantage of a metric, named kernel divergence. This metric has two

21  advantages: First, it directly compares dataset-wise distributional differences with no requirements on

22  human intervention or prior knowledge-based pre-classification. Second, minimal assumptions in

23  distribution and sample space are made in data processing to enhance its application domain. This enables

24  test-verified appropriate handling of datasets with both linear and non-linear relationships. The model was

25  then validated in a case study with eight SCRS phenotyping datasets each sampled from a different

26  enhanced biological phosphorus removal (EBPR) activated sludge community located across North

27  America. The model allows the determination of sufficient sampling size for any targeted or customized

28  information capture capacity or resolution level. For example, an approximated sampling size of 50 or

29  100 spectra for full-scale EBPR-related ecosystems at 5% or 2% OPU cluster resolution. Promised by its

30  flexibility and minimal restriction of input data types, the proposed method is expected to be a

31  standardized approach for sampling size optimization, enabling more comparable and reproducible

32  experiments and analysis on complex environmental samples. Finally, these advantages exhibit the

33  capability of generalizing to other single-cell technologies or environmental applications, provided that

34  the input datasets contain only continuous features.

35

36

37  **TOC**

38

39

40

41

42

43

44

45

46

## INTRODUCTION

Advances in various modern analytical methods such as single-cell technologies have enabled unprecedented high-resolution and fundamental study of environmental microbiology than traditional cultivation-based and bulk-measurement methods. Some examples include metabolite probing (e.g. stable isotope probing) [1,2], single-cell phenotype identification (e.g., fluorescence *in-situ* hybridization (FISH)) [3] and sensitive, high through-put cell sorting (e.g., FISH-activated flow cytometry (FACS) and optical tweezer-based cell sorting) [4,5]. Situational studies in addition exhibit demand for non-invasive, real-time, label-free and continuously observation methods, in complement to or beyond these current single-cell technologies. These technologies, including single-cell Raman microspectroscopy can reveal cell response and metabolic changes under stimulation from various environmental changes.

Being a member of vibrational spectroscopic technology, Raman spectroscopy profiles the photons which are inelastically scattered to different frequencies due to the energy exchange between the monochromatic photon and a vibrating molecule. Its result spectra encode the fingerprints for pinpointing the chemical composition in observed cell sample, and ultimately, resolving its cellular phenotype and metabolic state. Single-cell Raman spectroscopy (SCRS) and its combination with other single-cell methods present promises for meeting this demand [6-9], and they have been demonstrated as powerful techniques in sub-cellular level substrate composition profiling [10,11], high through-put metabolic pathway and cell type identification [6,12,13], cell sorting [14], and qualitative or quantitative 3D structural imaging [15,16].

SCRS has been explored and demonstrated as one of the top candidate single-cell techniques, for cell identification up to strain-level discrimination of microorganism members from targeted community from different environmental matrices, including clean room [17], drinks [18], food [19], water [20], or cerebrospinal fluid [21]. Xu and Webb et al. demonstrated that the high resolution of SCRS was able to discern two strains of only single-gene mutation apart [22]. In addition, in contrast and complementary to genomics-based microorganism profiling approaches, SCRS captures and reflects the cell's "metabolic state", which is

71    more dynamic and responsive to environmental stimuli. Furthermore, SCRS enabled technologies may

72    help fill the gap in linking cellular phenotypes with their genotypes [23].

73    One challenge associated with the application of single-cell technology such as SCRS for complex

74    environmental samples is the determination of sufficient sampling size without prior knowledge of the

75    system diversity, yet with the need to balance between statistical power and cost of resources and time.

76    This challenge raises from two major facts. First, to our knowledge, current automation level for SCRS

77    sampling still under-satisfies the high through-put requirements for very-large scale surveys. Restricted

78    by its labor and time demand, most previous studies randomly select a sample size (i.e., the number of

79    single-cell spectra to collect per environmental sample) [22, 24], estimate empirically [25], or follow lab-specific

80    protocols [11, 26]. A commonly selected range of SCRS dataset size is 200-1,000 in dependence to the

81    population complexity [23, 27, 28] or 20-200 spectra per label in classification-oriented studies [21, 25, 29] with a

82    largest reported total sample size of 10759 [30]. Second, the level of microbial community diversity in

83    different environmental samples varies largely, therefore it is difficult to estimate *a priori*. The optimal

84    SCRS sampling depth for any given system remains unsolved and it limits the standardization of SCRS

85    for its wider applications. This drives the demand for a robust method which statistically validates the

86    sampling depth without knowing the composition and complexity of the microbial community.

87    Sampling size assessment of SCRS from environmental samples were discussed in previous studies but

88    is often restricted to situational applications. Learning-curve (LC) based technique targeting 5% Bayes

89    error rate was proven effective to investigate proper sample size to train a classifier [31, 32]; however it is a

90    quite different objective and this method is not suitable for unsupervised applications. Majed et al. (2009)

91    first attempted a practical solution by iteratively sampling and classifying samples, tracking abundance

92    changes of classified categories [25]. Their relative abundances would be repeatedly calculated, adding a

93    fixed number of spectra each time, until they stabilized above a sample depth threshold. However, its

94    reliance on classification requires a significant amount of both human intervention and domain knowledge

95    to select appropriate discriminating criteria. For example, in their proposed protocol, an exploratory

96    experiment was first carried out to identify 65% biomass (in biovolume) as a functional group of known

97    metabolic traits. The choice of classification criteria which target this majority type of cells therefore was

98    validated. Despite the cost of preliminary experiments, such a dominating microbial group or species is

99    not guaranteed to exist in more complex environmental samples. introduction of an adopted rarefaction-

100    like technique for enabling direct application on continuous feature datasets (e.g., SCRS) without pre-

101    classification, which defined maximum pairwise Euclidean distance as the diversity measure of a sample

102    set, was named as "diversity index (DI)" [28]. However, two potential issues exist with this DI-based

103    approach. First, the DI of the entire observation dataset is directly treated as reference in sampling size

104    assessment, implicitly regarding it as the population but without further validation. Second, using

105    maximum Euclidean distance to represent the diversity discards all detailed distance distributional

106    structures, potentially causing under-estimation of the true diversity. Thirdly, the definition and parameter

107    for quantifying "diversity" is required, which may not be available such as the case for the SCRS data.

108    In this study, we propose a new algorithm for sample size assessment that circumvents the disadvantages

109    inherent to previously reported approaches. Our algorithm iteratively increases the number of samples

110    until the optimal sample size is achieved. At each increment, our algorithm measures the distributional

111    difference due to the increase in sample size via *kernel divergence*, a pseudometric that measures the

112    difference between two population distributions. As more samples are collected and observed, the

113    distributional change due to the added samples also converges towards zero. By observing this decreasing

114    trend, our algorithm is capable of predicting the sample size as the change in distribution becomes less

115    than a user-defined threshold. Compared to previously reported approaches, this method has two major

116    advantages. First, no distribution or linear relationship is assumed in the input data; this enables more

117    general and improved handling of datasets when such assumptions do not necessarily hold. Second, it is

118    unsupervised, granting the independence to either classification or other pre-processing which may

119    require extra prior knowledge (e.g., microbial community composition) and human intervention (e.g.,

120    selection of classification criteria). This method is validated in a case study using SCRS datasets sampled

121 from eight independent enhanced biological phosphorus removal (EBPR) microbial communities, each

122 obtained from a different North American wastewater treatment facility. To our knowledge, this is the first

123 population-blind sample size prediction and assessment method that has been applied on biospectroscopic

124 datasets. The outcome will facilitate wider, more standardized and more reliable application of advanced

125 analytical technologies such as SCRS for various environmental studies. In addition, since minimal

126 assumption is made with the input data, this approach can potentially be applied on sampling size

127 assessments with any other dataset as long as it only contains continuous features.

128 **METHODOLOGY**

129 Given a population $Y$, we wish to discover the smallest possible subset of the population that still

130 statistically represent the total population. If we denote the subset as $X$ it can be determined by increasing

131 the size of the subset until its internal statistics become stable, i.e., when adding more samples no longer

132 updates the distribution of the subset. Kernel Divergence ($\mathcal{D}_p$) allows one to measure the non-linear

133 distribution difference between two distributions. This study leverages this pseudometric to evaluate the

134 statistical change between increments of samples added to the subset. By iteratively tracking this value,

135 the size of the smallest sufficient subset can be determined when the distributional change becomes
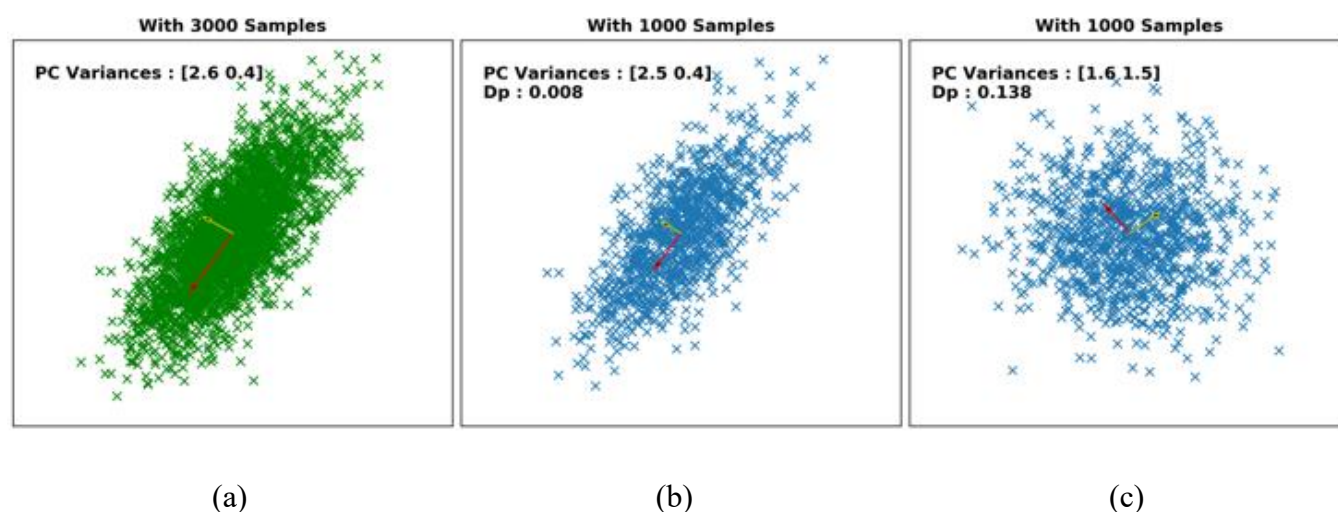
136 negligible.

137 **Kernel Divergence**

138 Inspired by He et al. (2017), we noticed that to reliably measure the distributional differences between

139 two sample sets is critical in sampling size assessment. There are currently many ways of measuring the

140 distributional difference between two populations. Broadly, they can be organized into two categories: F-

141 divergence [33] and Integral Probability Metrics (IPMs) [34]. F-divergence uses the ratio between two

142 distributions to measure their similarity while IPMs use their difference. In general, F-divergence requires

143 the researcher to know the distribution ahead of time. This makes it difficult to compute the divergence

144 given just samples, e.g. the Kullback-Leibler divergence [35]. Alternatively, IPMs do not require prior

145 knowledge of the sample distributions. Instead, they approximate the distribution directly from the

146 samples. A standard IPM is the Maximum Mean Discrepancy (MMD) . It measures the similarity between

147 two distributions by comparing their 1st moment in the Reproducing Kernel Hilbert Space (RKHS). In its

148 original space, it consequently compares all of its moments. Kernel divergence is different in that it uses

149 the 2nd moment in the RKHS, therefore, it is a shift and rotation invariant in RKHS. Instead of comparing

150 the distributions based on their moment, kernel divergence looks at the shape of the data in RKHS as

151 measured by its variance along its principal components. This approach is extremely memory efficient as

152 it allows us to compress the distributional information with a couple of eigenvalues. Any new incoming

153 population can consequently be also reduced to these numbers to compare their distribution difference.

154 We propose the usage of kernel divergence, motivated by two observations. First, by mapping data onto

155 the RKHS, the non-linear aspects of the data can be captured for analysis[36]. Second, the variances along

156 the principal components (PCs)[37] of a dataset summarize the shape of the data. **Figure 1** (a) and (b) shows

157 a scatterplot of two populations from the same distributions. Note that PC variances of the distributions

158 in (a) and (b) closely matches each other, indicating a resemblance in the shape of the data. Alternatively,

159     **Figure 1** (c) shows a scatter plot from a population belonging to a different distribution. Note that its

160     variances along the PCs are also noticeably different from that in (a) or (b).

161

          (a)                                    (b)                                    (c)

162

163    **Figure 1.** Example of how the principal components can be used to compare distributions. (a) and (b) were generated by an

164    identical population while (a) has three times the number of samples than (b). (c) was generated from a different distribution

165    and exhibited different variation ratios along PC1 and PC2 in comparison to (a) and (b).

166     Unfortunately, since the PCs can only capture linear relationships, it is no longer an appropriate tool when

167     the data fail to match the linear assumption. Since the data distribution from real applications are

168     commonly unknown, using PCs to compare distributions may not be appropriate for all cases. The idea

169     of Kernel Divergence is to combine RKHS's ability to capture non-linear relationship with PC's ability to

170     summarize the data. By first projecting the data into RKHS, non-linear PCs [38] can now be used to compare

171     the distributions. Since the concept of PCs is commonly used for the original data space, we will

172     distinguish the PCs in RKHS as the kernel principal components, or KPCs.

173     There are several advantages in using Kernel Divergence. First, the value computed by the divergence

174     can be treated as a distance between two distributions. It is always a positive value where a $\mathcal{D}_p = 0$

175     denotes a complete equivalence of the empirical variances along the KPCs. Conversely, a larger $\mathcal{D}_p$ also

176     indicates a further distance between two distributions. Second, the value of the divergence $\mathcal{D}_p$ has

177     practical meanings as suggested by its theoretical proof. Namely, it denotes the worst-case error along a

178    KPC. For example, if $\mathcal{D}_p = 0.01$, then the biggest difference between two samples along any single KPC

179    is bounded within 1%. Third, and more importantly to practitioners, $\mathcal{D}_p$ can be efficiently computed using

180    only a few lines of code with existing open-source software.

**Computing the Kernel Divergence**

182    We define a population of samples as $X \in \mathbb{R}^{n \times d}$ where $n$ and $d$ denotes the number of samples and the

183    dimension respectively. Let $H$ be a centering matrix defined as $H = I_n - \frac{1}{n}\mathbf{1}_{n \times n}$ where $I_n$ denotes an

184    identity matrix of size $n$ and $\mathbf{1}_{n \times n} \in \mathbb{R}^{n \times n}$ denotes a matrix of 1s. Then a centered kernel matrix is

185    defined as $HK_XH$ where i-th row and j-th column element of the $K_X$ is defined as

$$K_{X_{i,j}} = \exp\left\{ -\frac{\|x_i - x_j\|_2^2}{2\sigma^2} \right\},$$

187    known as radial basis function (RBF) kernel, where $\exp(\cdot)$ stands for the exponential function with base

188    $e$, and $\|x_i - x_j\|_2^2$ is the squared Euclidean distance between i-th and j-th samples. We chose this kernel

189    function for its flexibility to approximate a wide range of non-linear functions.

190    Given the definition of a centered kernel, we define $K_\mathbb{A}$ and $K_\mathbb{S}$ as the centered kernels for two population

191    samples of $\mathbb{A}$ and $\mathbb{S}$. Let the $m$ largest eigenvalues of $K_\mathbb{A}$ be $\lambda_\mathbb{A} = \left[ \lambda_\mathbb{A}^{(1)}, \dots, \lambda_\mathbb{A}^{(m)} \right]$, where $\lambda_\mathbb{A}^{(1)} \geq \cdots \geq$

192    $\lambda_\mathbb{A}^{(m)}$, and the $m$ largest eigenvalues of $K_\mathbb{S}$ be $\lambda_\mathbb{S} = \left[ \lambda_\mathbb{S}^{(1)}, \dots, \lambda_\mathbb{S}^{(m)} \right]$, where $\lambda_\mathbb{S}^{(1)} \geq \cdots \geq \lambda_\mathbb{S}^{(m)}$. Here, $m$ is

193    preferably the number of "major" eigenvalues indicated being before a significant value drop in the

194    eigenvalue spectra plot. Eigenvalues are the variances in the respective eigenvector directions. It is notable

195    that some datasets will exhibit gradually decreasing eigenvalues with no such drop. In such cases, an extra

196 parameter $p$, the percentage of the total variance which we wish to preserve from the population, is

197 introduced to calculate $m$; specifically,

$$198 \qquad m = \min\left\{a \mid \frac{\sum_{i=1}^{a} \lambda_{\mathbb{A},i}}{\sum_{i=1}^{a} \lambda_{\mathbb{S},i}} \geq p\right\}.$$

199 Once the centered kernels of the two populations are calculated, the kernel divergence between $K_{\mathbb{A}}$ and

200 $K_{\mathbb{S}}$ is defined as

$$201 \qquad \mathcal{D}_p(K_{\mathbb{A}}, K_{\mathbb{S}}) = \left\| \frac{1}{\|\lambda_{\mathbb{A}}\|_1} \lambda_{\mathbb{A}} - \frac{1}{\|\lambda_{\mathbb{S}}\|_1} \lambda_{\mathbb{S}} \right\|_{\infty},$$

202 where $\|\cdot\|_1$ and $\|\cdot\|_\infty$ are $L_1$ and $L_\infty$ norms respectively. A brief proof of kernel divergence as a measure

203 is provided in Supporting Information **Proof S1**.

204 **Sample Size Assessment**

205 Assume that a sample set $\mathbb{S}$ has already been physically acquired from a population $\mathbb{P}$. However, knowing

206 only this end-status of sampling process but no *a priori* population composition knowledge for reference,

207 sample size assessment would be almost impossible. To overcome this difficulty, we model the sample

208 acquisition process for better profiling of the sample distribution changes via randomized virtual sampling

209 simulation; while, this modelling step implicitly assumes that the physical sampling experiment (which

210 acquires $\mathbb{S}$) is unbiased. This will result in a kernel divergence profile at each sampling depth, during

211 acquisition towards the same dataset $\mathbb{S}$ while in a randomized order. Such simulation can be repeated

212 multiple times for estimation of the mean kernel divergence profile with minimal dependence to the

213    randomization effects. These modelling steps summarize the core algorithm in our assessment protocol,

214    represented as in the pseudo-code below:

215      (1) $A_1 \leftarrow$ randomly select $k_0$ samples from $\mathbb{S}$;

216      (2) $A_{n+1} \leftarrow$ randomly select $k$ samples from $\mathbb{S} \backslash A_n$, insert to $A_n$;

217      (3) Calculate kernel divergence $\mathcal{D}_p(A_n, A_{n+1})$;

218      (4) Repeat (2)-(3) until $\mathbb{S} \backslash A_n = \emptyset$;

219      (5) Repeat (1)-(4) multiple times; this results in multiple kernel divergence profiles by acquiring the

220         same dataset $\mathbb{S}$ in different orders.

221    For simplicity, at each Step (2), the number of samples drawn is standardized as a fixed number $k$, denoted

222    as *batch size*. $k$ is often determined by the physical analytical system. This is a sensitive parameter to

223    kernel divergence calculation, thus should be determined *a priori* and kept unchanged during a single

224    assessment. For SCRS data, since the Raman system yields single spectrum for each cell at each sampling

225    event, the $k$ value is therefore 1. Finally, we determine sample sufficiency of dataset $\mathbb{S}$ by checking if the

226    "average" kernel divergence profile resulted from the above steps has converged to zero with a pre-defined

227    threshold $t$. Applying the interpretation of kernel divergence, the presence of such point of convergence

228    (POC) identifies a sampling depth at which further addition of $k$ more samples (i.e., a *batch*) will no

229    longer significantly update the sample distribution. In other words, the presence of such POC implies that

230    $\mathbb{S}$ is *sufficient*, otherwise *not sufficient*. The eigenvalue preserving percentage $p$ (if used in kernel

231    divergence calculation) and the number of iterations in Step (5) (default: 1000) are the last two adjustable

232    parameters in our proposed protocol, in addition to $k$ and $t$.

233    **SCRS Datasets from EBPR Facilities**

234    SCRS-based phenotypic survey was conducted on 8 different EBPR-related sludge samples, each

235    generating an individual SCRS dataset. Investigating the community composition phenotypic

236    characteristics captured by Raman spectra is conceptually analogous to the operational taxonomic units

237 (OTUs) based survey via 16S-rRNA amplicon sequencing [23]. Each sludge sample represents the microbial

238 community in the EBPR anaerobic reactor of a different wastewater reuse and reclamation facility

239 (WRRF) across the North America, with various geological, configurations, operational and influent

240 water characteristics (**Table S2**), including 4 conventional EBPR and 4 side-stream EBPR processes.

241 Studies showed that Raman spectra are sensitive to the experimental conditions and instrumental factors

242 therefore all an identical protocol was followed in acquisition of those 8 SCRS datasets to maximize their

243 cross-comparability. Briefly, each sludge sample was independently performed a phosphorus release and

244 uptake kinetics batch test as described by Gu et al. (2008) [39]. Raman-based phenotypic survey was

245 performed on the sludge extracted throughout the batch test following the preparing and acquisition

246 protocol described by Majed et al. (2008) and Onnis-Hayden et al. (2019) [40-42]. All spectra were acquired

247 with a 400-1800 cm$^{-1}$ range which is often referred as the "fingerprint range" for various cellular or

248 biomass substances [6, 23, 43]. All acquired spectra were then preprocessed with LabSpec 6 (HORIBA, 2

249 Miyanohigashi, Kisshoin, Minamiku Kyoto 601-8510 Japan), for cosmic spike removal, smoothing,

250 background subtraction, baseline correction and vector-normalization. The survey resulted in a total of

251 922 spectra in eight WRRF-specific datasets labelled from A-H. Two datasets, F (207 spectra from

252 Westside Regional, S2EBPR) and H (214 spectra from Upper Blackstone, conventional EBPR) had larger

253 sample size in comparison to the other six (ranging from 80-89).

254

255 **RESULTS AND DISCUSSION**

256 The proposed protocol was first tested with synthetic datasets to demonstrate its sample size assessment

257 performance and results interpretation. Multiple tests with different parameter settings were conducted

258 for discussion on the selection of parameters with respect to the experimental protocol in real applications.

259    Then the model was validated in a case studying with microbial SCRS phenotyping datasets acquired in

260    experiments.

261    **Test with Synthetic Datasets**

262    The tested synthetic datasets include both simple, linear and non-linear relationships. The dataset

263    construction and the test results are as follows:

264    ***Dataset construction.*** The first dataset (Dateset 1) is a series 2-dimensional datasets to test the model's

265    performance on sample size assessment (**Figure 2** (left)). A total of five sub-datasets were generated from

266    an identical mixture distribution to resemble five independent sampling experiments from a same, infinite-

267    sized population but targeting at different sampling sizes, respectively 20, 40, 100, 200 and 500. The

268    source distribution is composited by four 2-dimensional Gaussian distributions (i.e. $D = 2$, $K = 4$), and

269    is designed to be non-overlapping and linearly separable for better testing our model performance with

270    simple datasets. Therefore, each Gaussian sub-population is centered respectively at ($\pm$2.5, $\pm$2.5), with

271    variance of 0.1 on both axes.

272    Dataset 2 has a spiral morphology [44] aimed for testing with complex, non-linear datasets. This 2-

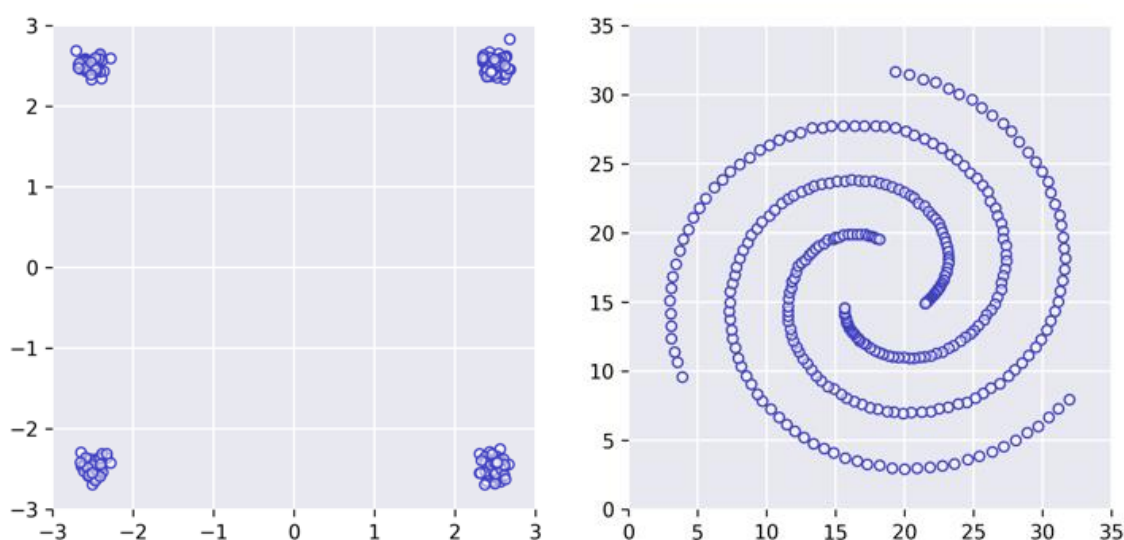273    dimensional dataset contains 3 clusters in total of 312 samples ($D = 2, N = 312, K = 3$) (**Figure 2**

274    (right)).

275



276    **Figure 2.** Visualization of synthetic datasets used in performance testing of proposed sample size assessment method using

277    kernel divergence. Dataset 1 (left) is generated with a mixture composed of four 2-dimensional Gaussian sub-populations, each

278    has 25% of the total population and independent variance of 0.1 on both axes. It contains five sub-datasets of various sizes (20,

279    40, 100, 200 and 500). The 100-sample dataset is shown as a representative. Dataset 2 "spiral" (right) is a 312-sample, 2-

280    dimensional dataset containing three clusters, each forms a non-linear spiral shape [44].

281    ***Sufficient sampling size assessment.*** We first demonstrate the calculation and decision making of the

282    proposed sampling size assessment protocol with Dataset 1. The kernel divergence at different sampling

283    depths are shown in **Figure 3**, calculated using empirical parameter settings with batch size $k = 1$ to

284    simulate the one-by-one sampling strategy, and using empirical convergence threshold $t = 0.01$ and

285    1,000 iterations. We determined $m = 3$ as almost 100% variation fell on the first 3 KPCs (**Figure S2**).

286    The selection of threshold $t = 0.01$ was identified in sensitivity test corresponding to a 0.8% foreign class

287    abundance (will be discussed in detail later). The mean kernel divergence was plotted as solid lines shaded

288    with the standard deviations observed at each sampling depth.

289    The closely overlapping kernel divergence profiles indicates that POC estimation is rather robust and

290    consistent, being independent to the dataset size but only the sample distributions. For this data set, the

291    estimated sufficient samples size is around 40 with a targeted convergence threshold of 0.01. As shown

292    in Figure 3, sample size less than 40 ( N=20, 40) was not sufficient to meet the targeted POC. Of course,

293    the sufficient sample size depends on the user-chosen POC level. If the target POC threshold is at 0.005,

294    then the sufficient samples size for this data set would be around 85.
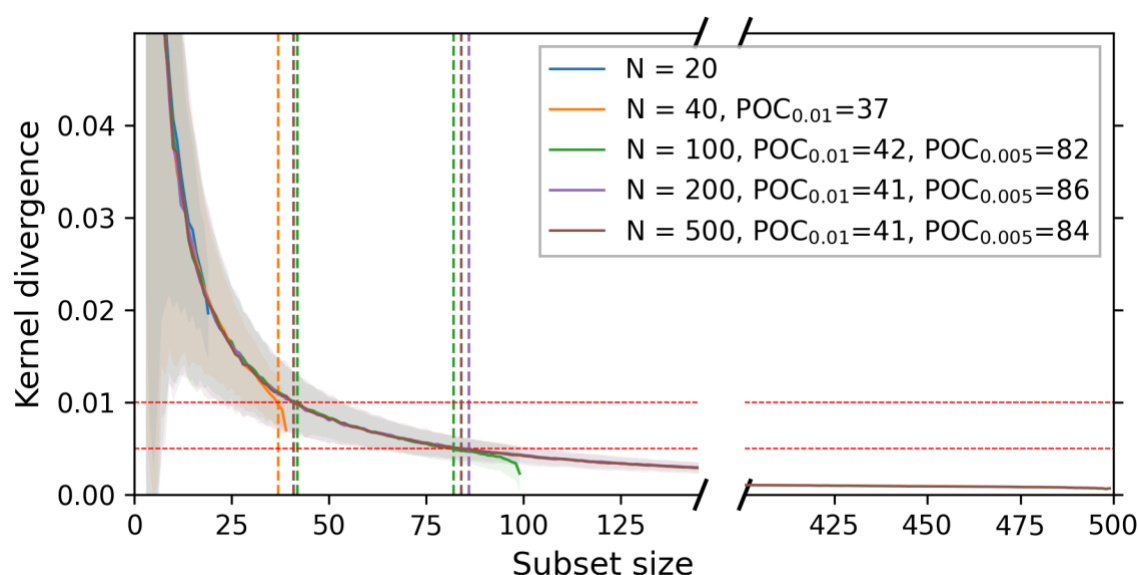
295



296

**Figure 3.** Per sample kernel divergence profiled at various sampling depths simulated independently from each sub-dataset in

Dataset 1. All simulations used first 3 eigenvalues $m = 3$, batch size $k = 1$, 1000 iterations and convergence thresholds $t =$

$0.01, 0.005$ (shown as horizontal lines). The point of convergence (POC) indicates the minimal sample size required to bound

the distributional difference below the indicated threshold, subject to adding one more sample. A lower threshold corresponds

to larger sample size.

302    ***Assessment with non-linear dataset.*** A series of similar simulations were conducted on the Dataset 2

303    "spiral" containing subpopulations that are not linearly separable. All assessment simulations were

304    conducted with 1,000 iterations and three different batch sizes $k = 1, 2, 5$. We chose $m = 5$ eigenvalues

305    which encode approximately 95.5% of total variances along KPCs (**Figure S3**) . A convergence threshold

306    $t = 0.005$ was selected for this case, which was identified in sensitivity test corresponding to a 3.3%

307    foreign class abundance. Results are shown in **Figure 4**. The sufficient sample size varied depending on

308    the batch sample size k values. Note that k value is physically determined by the analytical method itself.

309    Here, we demonstrated that the sufficient samples size for targeted POC threshold increased from 59 to

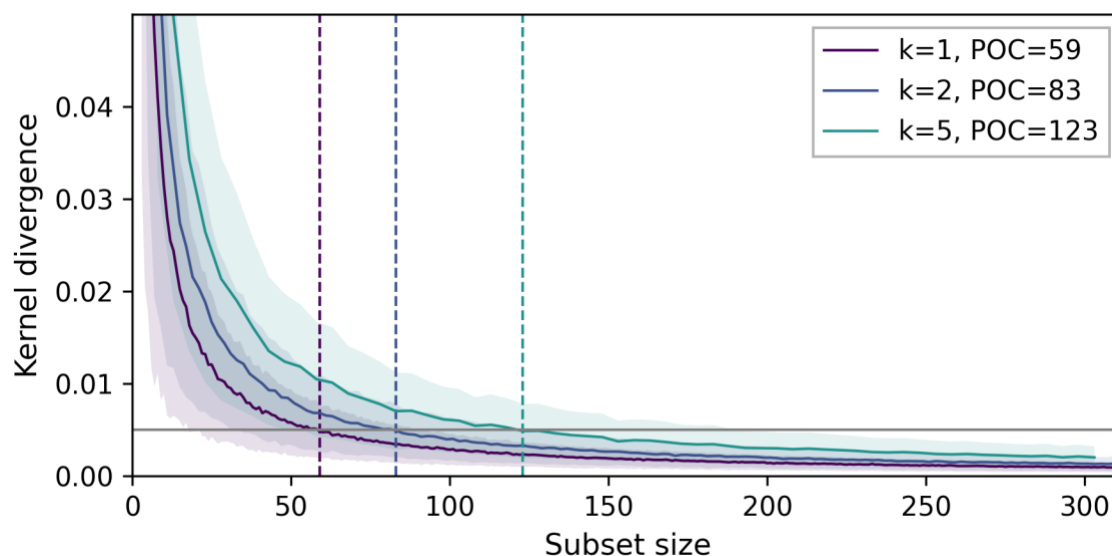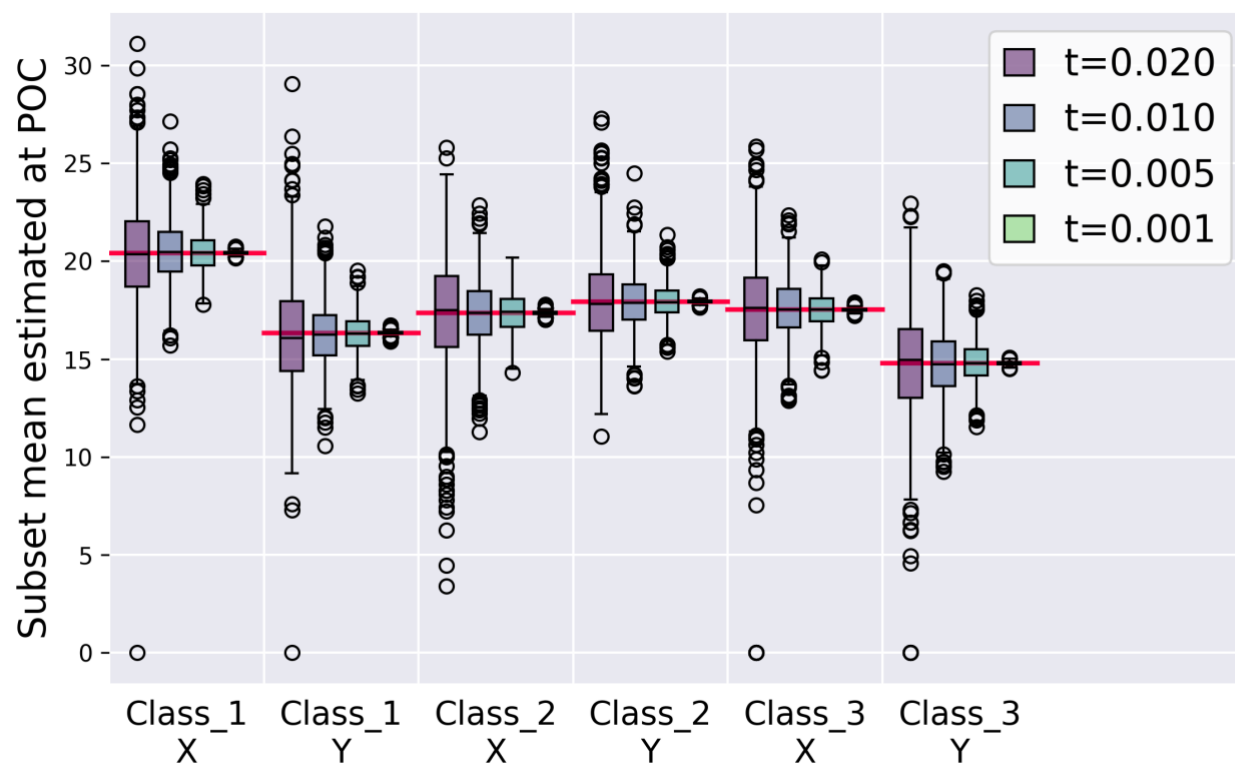310    123 as the batch sample size increased from k=1 to k=3.



311

312    **Figure 4.** Sample size assessment results with non-linear dataset. Plot shows kernel divergence profiled with different number

313    of samples (1, 2 and 5) per iterative addition (batch size), simulated with Dataset 2 "spiral". All calculation used first 5

314    eigenvalues ($m = 5$) and were from 1000 simulation iterations. Convergence threshold $t = 0.005$ (shown as horizontal line).

315    The point of convergence (POC) indicates the minimal sample size required to bound the distributional difference below the

316    indicated threshold, subject to each batch size.

317    **Figure 5** investigates the statistical properties of the POCs in each simulation iterations under different

318    convergence thresholds, which were identified as the last step in each specific simulation iteration having

319    a kernel divergence larger than the given threshold. These results showed that, taking the advantage of

320    RBF kernel, our method can effectively capture the distributional information in a dataset with non-linear

321    aspects. For example, with convergence threshold $t = 0.005$, 2.5%-97.5% quantile of identified POC

322    subsets had estimated class means being within $\pm2.03$ of the actual population mean on both dimensions.

323    And, the estimated class abundances were within $\pm 7.0\%$ of actual abundance. At 25%-75% quartile, error

324    of class means, and abundances were $\pm 0.72$ and $\pm 2.7\%$ respectively. A smaller threshold would result in

325    more parametrically accurate POC subsets; however, it requires more samples to achieve reliable POCs.

326    The test showed that the POC subset sizes had 25%-75% quartile respectively from 47-61 with $t = 0.01$,

327    102-124 with $t = 0.005$ or 307-312 with $t = 0.001$. In addition, incorporating the radial basis function

328    (RBF) kernel enabled the kernel divergence to be appropriate in measuring dissimilarity between two

329    datasets with non-linear aspects, and therefore our method is generally applicable in other datasets.
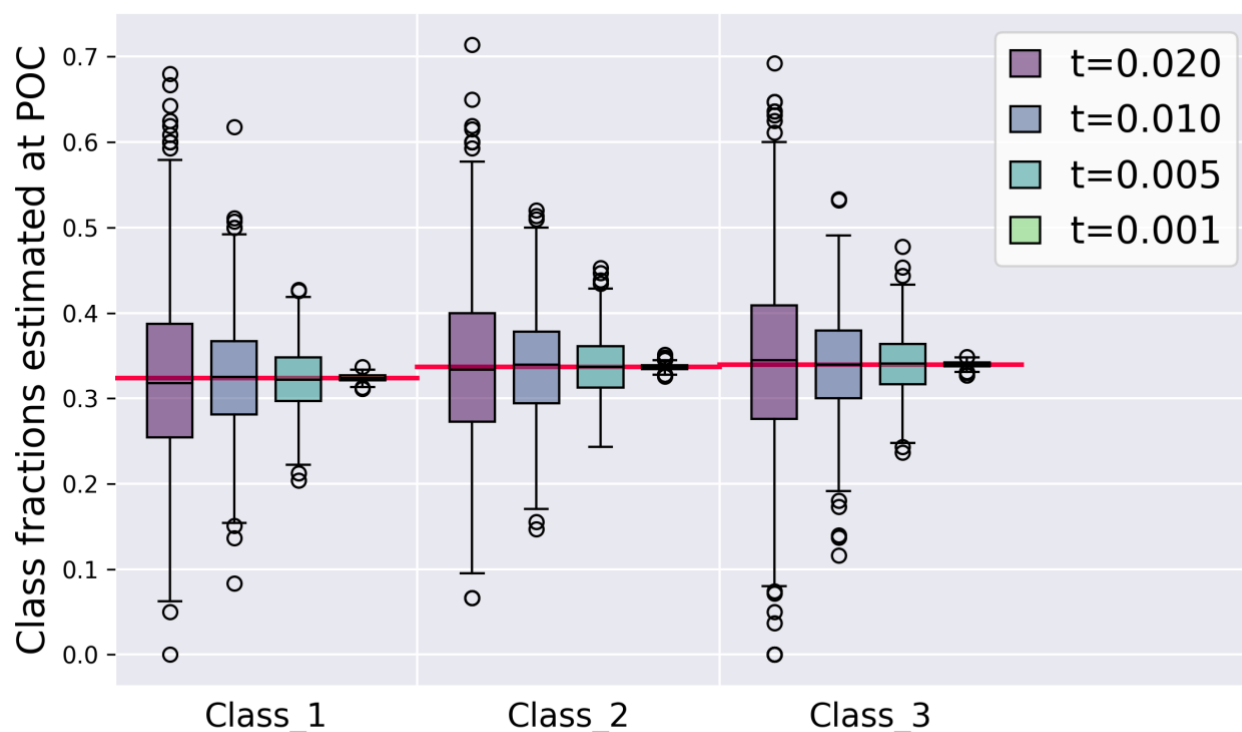


330

331

**Figure 5.** Box plots of mean (top) and fractions (bottom) for each class in the subset selection at each individual POC with different convergence threshold $t$. from the 1,000 simulation iterations. The simulation POC is determined as the last step in specific sampling simulation that had a kernel divergence larger than the given threshold $t$. The horizontal red lines indicate the ground truth calculated from the whole "spiral" dataset.

***Selection of convergence threshold t.*** Results in **Figure 5** showed the relationship and potential impact to the mean and abundance estimations of each class associated with different convergence threshold criteria. As the threshold pre-determines the resolution of two datasets being asserted "different", it also relates with the minimal abundance of identifiable classes. To reveal this effect, we chose one class from the original datasets, Dataset 1 and 2, (referred to as "the foreign class"), and randomly delete a sub-selection of its samples. This creates artificial datasets with the foreign class at varying abundances. With random repeats and the dataset with no foreign class as a reference, we could estimate the minimal abundance of the foreign class to exhibit a significant change in kernel divergence under chosen threshold. Note the choice of foreign class has minimal impact to the results due to the high symmetry. The results

345    in **Figure 6** showed that under a same kernel divergence threshold level, the size of detectable foreign

346    class varies largely depending on the parent dataset. The simpler dataset (Dataset 1) requires less samples

347    to achieve the same kernel divergence level comparing to the non-linear Dataset 2. For example, mixing

348    foreign class at abundances of 0.4% and 0.8% leads to 0.005 and 0.01 kernel divergence; while 3% and

349    9.5% abundances of the foreign class would be required respectively with Dataset 2. Considering both

350    the effect of precision, resolution and dataset complexity, we suggest that $t \leq 0.01$ would be adequate

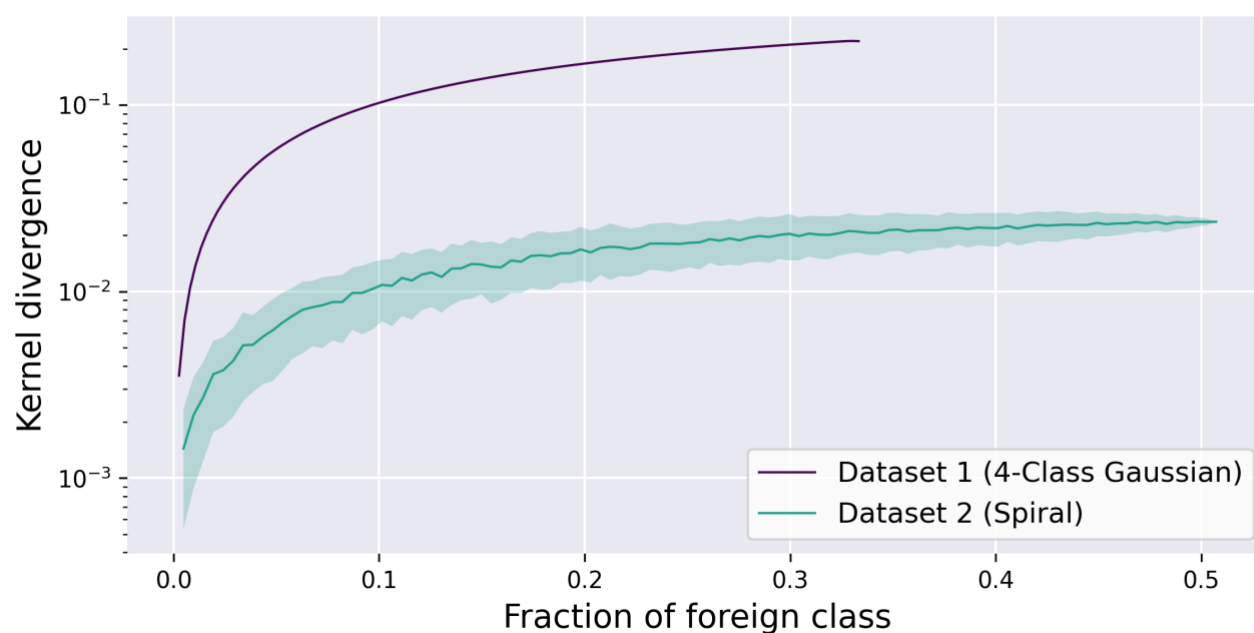351    empirically in general, and $t \leq 0.001$ is considerable when a high accuracy is desired.



352

353    **Figure 6.** Kernel divergence sensitivity test with varying the abundance of one of the classes (referred to as "foreign class").

354    All kernel divergences were calculated in reference to the absence of the entire foreign class. The lines and shades show the

355    mean and standard deviation respectively estimated by random sub-selection of samples in the foreign class. The standard

356    deviation in the Dataset 1 profile is minimal and barely visible.

357    **Case Study: Sample Size Assessment for SCRS Phenotyping Datasets**

358    After performance evaluation on two synthetic datasets with or without linear relationship, the proposed

359    algorithm and method were applied to investigate sampling size requirements with single-cell Raman

360    spectroscopic dataset retrieved from eight microbial communities representing eight different wastewater
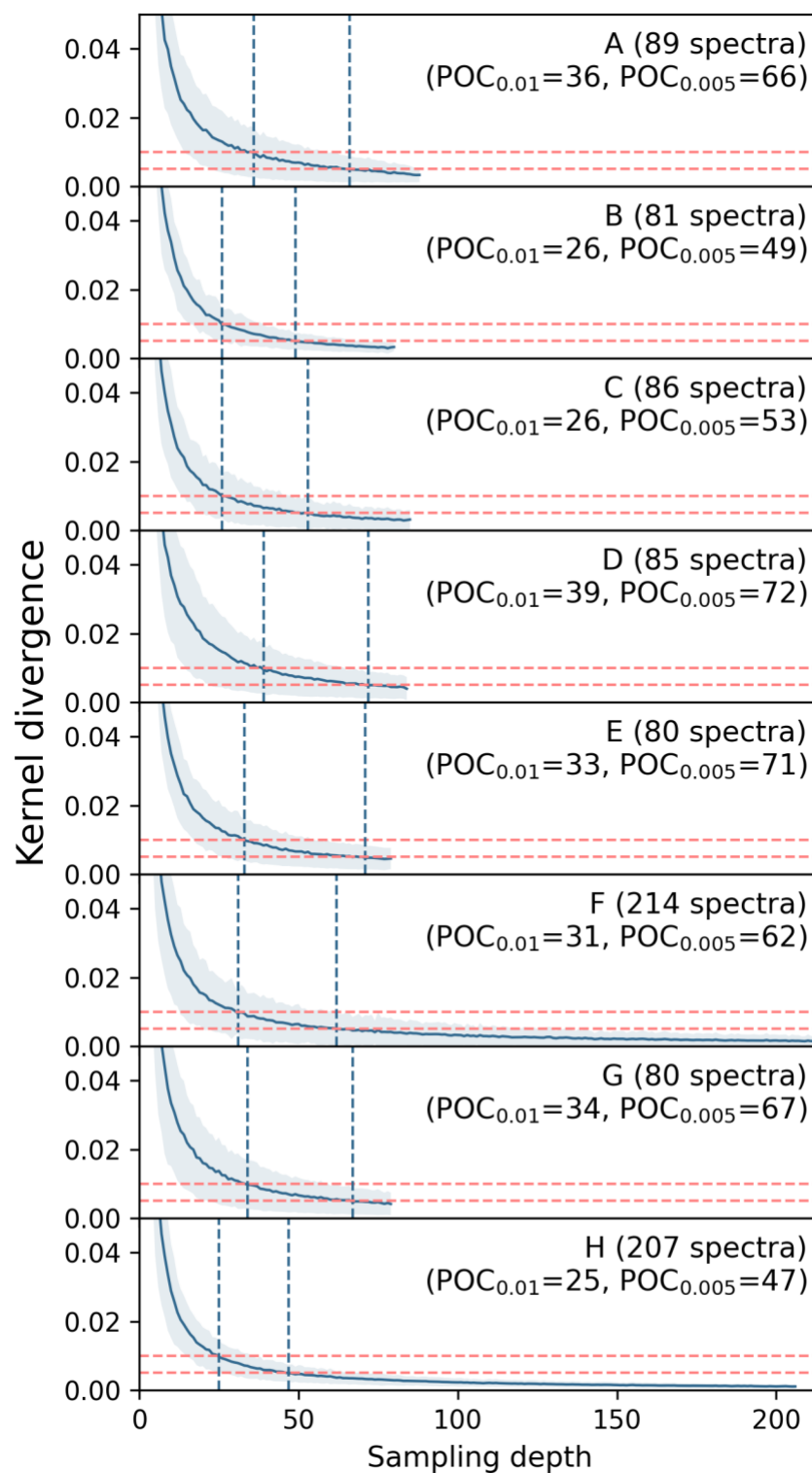
361    reuse and reclamation facilities (WRRFs) located in North America [27, 42]. Details on these EBPR facilities,

362    sampling and SCRS data acquisition were described in the methods section and in supporting information.

363    ***Sampling size assessment.*** Kernel divergence profiling simulations were conducted independently on

364    individual SCRS-based phenotyping dataset, using parameter $k = 1$ , $t = 0.01, 0.005$ and 1000

365    permutations. Sample batch size k=1 since our data was obtained using single-cell Raman

366    microspectroscopy at single cell resolution [27]. As no clear group of major eigenvalues can be identified,

367    we used parameters $p = 96\%$ in calculations (**Figure S4**).

368    As shown in **Figure 7**, the sufficient sample size based on each individual kernel divergence profile for

369    the eight EBPR communities ranged from 26-39 under a resolution threshold of $t = 0.01$. The parallel

370    assessment under a higher resolution by lowering the threshold of convergence to $t = 0.005$ revealed an

371    increased sample size range from 47-71 among the EBPR plants. These were below the empirical

372    reliability checking criteria discussed previously (POC: $N < 0.85$). This range are consistent with

373    previously reported values (60-65 samples) proposed by Majed et al. (2009) via an alternative rationale

374    [25].

375    ***Investigation of the convergence threshold.*** The sensitivity test shown in **Figure 8** further investigated

376    the convergence thresholds towards a more physical and practical interpretation, by evaluating the

377    maximum abundance level of sample clusters that could be potentially missed with varying thresholds.

378    Operational phenotypic units (OPU) clustering was first carried out individually to identify cluster groups

379    in each dataset, using correlation distance, average linkage as described by Li et al. (2018) [23]. **Figure 8**

380    then shows the maximum kernel divergence observed when randomly removing one identified OPU

381    cluster which is below a targeted abundance level. The results indicate that for 7 out of the 8 datasets,

382    using a threshold of 0.01 captures all OPU clusters of more than 5% abundance; and using a threshold of

383    0.005, this resolution can be improved to 2%. Therefore, we conclude that for these 8 datasets from EBPR

384    communities, the sample sizes were sufficient to capture all OPU clusters of at least 5% relative
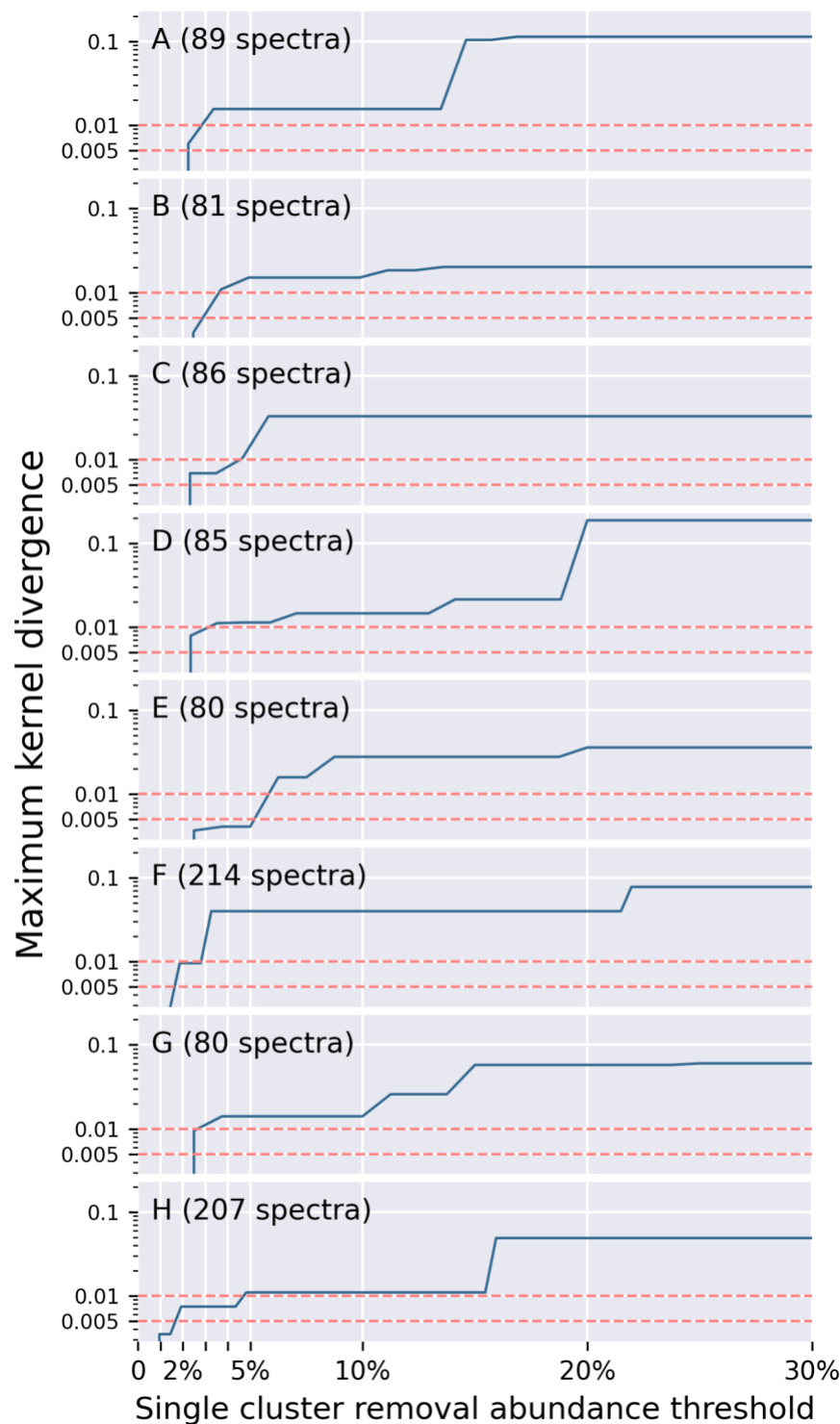
385    abundance, and among them, 5 date sets were further sufficient at capturing OPUs with 2% relative

386    abundance.



387

388

389    **Figure 7.** Kernel divergence profiles versus sampling depth, simulated independently on 8 single-cell Raman spectroscopic

390    (SCRS) microbial phenotyping datasets, from 8 individual full-scale enhanced biological phosphorus removal (EBPR) systems

391    in different wastewater reuse and reclamation facilities (WRRFs) located across North America. Simulation parameters were

392    eigenvalue preserving percentage $p = 96\%$, 1000 iterations, convergence threshold $t = 0.01, 0.005$ (shown as the two

393    horizontal lines) and batch size $k = 1$ for single-cell Raman microspectroscopic method. The point of convergence (POC)

394    indicates the minimal sample size required based on the targeted convergence threshold $t$ and batch size $k$.

**Figure 8.** Sensitivity test of kernel divergence convergence over single OPU cluster abundance. The kernel divergence values were calculated by randomly removing a single OPU cluster not exceeding an abundance threshold (x-axis). The maximum kernel divergence calculated at each abundance threshold was shown on the y-axis. The OPU clusters were identified using correlation distance and average linkage as described by Li et al. (2018) [23] using a same cut-off at 0.7. The results indicated

400    the relationship between OPU resolution and choice of convergence threshold $t$. For example, a point x=5%, y=0.01 means

401    that using a threshold $t = 0.01$ allows at most one OPU up to 5% abundance being ignored. Such assessment results may be

402    specific to each dataset.

403    It is noticed that the identified sample sizes for the same target POC threshold among the eight EBPR

404    communities were rather comparable, indicating an intrinsic "similarity" of the microbial phenotypic

405    "richness" in these full-scale EBPR systems in North America. The correlation between SCRS-based

406    phenotypic clusters (i.e. OPUs) and their phylogenetic OTUs, with underlying implications of the

407    discriminative power of Raman spectrum features for discerning cells at various taxonomic levels (i.e.

408    species, strains etc.), is still under investigation. Promising cell identification at strain level have been

409    reported (ref.). The comparable kernel divergence profiles and narrow range of minimal sample size for

410    a given targeted POC threshold for the 8 EBPR microbial communities suggested that these engineered

411    wastewater treatment systems may have similar microbial phenotypic diversity measurements. How the

412    phenotyping profiles correspond to their phylogenetic composition are yet to be revealed and is beyond

413    the focus and scope of this study.

414    We also compared our results with another prior-knowledge independent method proposed previously by

415    He et al. (2017) [28]. The diversity measure of a given sample set, named as "diversity index (DI)" resides

416    as the core concepts in He's assessment protocol, which was defined as the maximum pairwise Euclidean

417    distance within the sample set. Through repeated virtual sampling experiments in a similar process, the

418    "average DI" at each sampling depth was estimated, then plotted as shown in supplementary **Figure S6**.

419    Finally, two sample size guidelines can be determined according to the decision criteria proposed by He

420    et al. (2017) [28]:

421    (1) 9-15 samples as "minimal size" identified when DI change per increasing sample depth by one is

422        less than 0.01 of maximal (i.e. dataset-wise) DI;

423    (2) 9-25 samples as "safe size" identified when average subset DI reaches 90% of maximal (dataset-

424        wise) DI.

425    These results were significantly smaller than the values identified by our protocol, indicating DI led to a

426    less diverse estimation in comparison to kernel divergence. Two potential reasons may have contributed

427    to different performances of the DI-based and kernel divergence-based protocols. First, DI ignores details

428    of the sample distance distribution but only its maximum, while kernel divergence utilizes comprehensive

429    information of the entire distance matrix. Therefore, DI-based calculations will probably make false

430    conclusions when two sample sets have different sample distributions but rather similar DI values.

431    Second, a sole reliance on the maximum distance also increases its sensitivity to the presence of outliers;

432    therefore, appropriate and sophisticated outlier removal techniques might also be necessary in real

433    applications. Comparing to kernel divergence-based protocol, the DI-based method was likely

434    underestimating the true diversity and complexity in our SCRS phenotyping datasets, resulting in reduced

435    reliability.

436    One of the main challenges in wider application of new emerging high-resolution technologies for

437    profiling and characterization of complex environmental systems, such as SCRS and cell imaging, is the

438    standardization of the experimental protocols and data analysis such as the optimal sampling size with the

439    consideration of both time and resources cost and information sufficiency. There is no widely accepted

440    approach and method for determining the sufficient sampling size on environmental datasets without pre-

441    classification. We proposed and validated a sample size assessment protocol using kernel divergence, a

442    novel dissimilarity measure at the dataset-level, which is a more comprehensive and systematic

443    quantitative comparison between two observation datasets. More importantly, our proposed method

444    enables the decision on sampling size without prior knowledge of the diversity and complexity of the

445    system. This property is especially powerful as demanded by *de novo* studies with environmental samples.

446    In addition, our proposed method has no restrictions on the input data as long as it contains continuous

447    features. In particular, it can capture data with linear and non-linear relationships. All these generalities

448    profit expansion to further potential applications. First, it provides a universal standard to compare the

449    sampling size determining criteria among different experiments in different labs, contributing to more

450    reliable, comparable and reproducible studies using similar single-cell technologies. In addition, robust

451    cross-comparison among different experimental protocols could be validated as well. Second, we believe

452    that the proposed sampling size assessment approach can be easily generalized to dataset generated from

453    other analytical technologies . Potential examples include Raman-based spectral histopathological

454    assessments, validating gating strategies in flow cytometry and collecting comprehensive cellular imaging

455    library based on visual or morphological measurements.

456    **Acknowledgement**

461

462    **ASSOCIATED CONTENT**

463    **Supporting Information**

464    **Proof S1.** Proof for kernel divergence properties.

465    **Figure S2.** Eigenvalue spectra plot for the 4-cluster Gaussian dataset (synthetic Dateset 1).

466    **Figure S3.** Eigenvalue spectra plot for the "spiral" dataset (synthetic Dateset 2d).

467    **Figure S4.** Eigenvalue spectra plot for the Upper Blackstone and Westside Regional datasets.

468　　**Table S5.** Supplementary information of WRRFs from which the SCRS datasets were sampled.

469　　**Figure S6.** Comparative analysis to the minimal sampling depth and safe sampling depth simulated

470　　following the approach proposed by He et al. (2017).

471

472　　**AUTHOR INFORMATION**

473　　**Corresponding Author**

474　　* aprilgu@cornell.edu

475　　**Author Contributions**

476　　‡These authors contributed equally.

477　　**REFERENCES**

478　　1.　　Radajewski, S., et al., Stable-isotope probing as a tool in microbial ecology. Nature, 2000.

479　　403(6770): p. 646-649.

480　　2.　　D Wang, P He, Z Wang, G Li, N Majed, AZ Gu. "Advances in single cell Raman spectroscopy

481　　technologies for biological and environmental applications." Current Opinion in Biotechnology, 64: p.

482　　218-229.

483　　3.　　Pernthaler, A., J. Pernthaler, and R. Amann, Fluorescence in situ hybridization and catalyzed

484　　reporter deposition for the identification of marine bacteria. Appl. Environ. Microbiol., 2002. 68(6): p.

485　　3094-3101.

486　　4.　　Kalyuzhnaya, M.G., et al., Fluorescence in situ hybridization-flow cytometry-cell sorting-based

487　　method for separation and enrichment of type I and type II methanotroph populations. Appl. Environ.

488　　Microbiol., 2006. 72(6): p. 4293-4301.

489　　5.　　Huang, W.E., A.D. Ward, and A.S. Whiteley, Raman tweezers sorting of single microbial cells.

490　　Environmental microbiology reports, 2009. 1(1): p. 44-49.

491    6.    Fernando, E.Y., et al., Resolving the individual contribution of key microbial populations to

492    enhanced biological phosphorus removal with Raman–FISH. The ISME journal, 2019. 13(8): p. 1933-

493    1946.

494    7.    Dina, N., et al., Rapid single-cell detection and identification of pathogens by using surface-

495    enhanced Raman spectroscopy. Analyst, 2017. 142(10): p. 1782-1789.

496    8.    Butler, H.J., et al., Using Raman spectroscopy to characterize biological materials. Nature

497    protocols, 2016. 11(4): p. 664.

498    9.    Harz, M., P. Rösch, and J. Popp, Vibrational spectroscopy—A powerful tool for the rapid

499    identification of microbial cells at the single‐cell level. Cytometry Part A: The Journal of the

500    International Society for Analytical Cytology, 2009. 75(2): p. 104-113.

501    10.   Moudříková, S.a.r., et al., Quantification of polyphosphate in microalgae by Raman microscopy

502    and by a reference enzymatic assay. Analytical chemistry, 2017. 89(22): p. 12006-12013.

503    11.   Wang, T., et al., Quantitative dynamics of triacylglycerol accumulation in microalgae populations

504    at single-cell resolution revealed by Raman microspectroscopy. Biotechnology for biofuels, 2014. 7(1):

505    p. 58.

506    12.   Lorenz, B., et al., Cultivation-free Raman spectroscopic investigations of bacteria. Trends in

507    microbiology, 2017. 25(5): p. 413-424.

508    13.   Ando, J., et al., High-speed Raman imaging of cellular processes. Current opinion in chemical

509    biology, 2016. 33: p. 16-24.

510    14.   Song, Y., H. Yin, and W.E. Huang, Raman activated cell sorting. Current opinion in chemical

511    biology, 2016. 33: p. 1-8.

512    15.   Kallepitis, C., et al., Quantitative volumetric Raman imaging of three dimensional cell cultures.

513    Nature communications, 2017. 8(1): p. 1-9.

514    16.   Freudiger, C.W., et al., Label-free biomedical imaging with high sensitivity by stimulated Raman

515    scattering microscopy. Science, 2008. 322(5909): p. 1857-1861.

516     17.    Rösch, P., et al., Chemotaxonomic identification of single bacteria by micro-Raman spectroscopy:

517     application to clean-room-relevant biological contaminations. Appl. Environ. Microbiol., 2005. 71(3): p.

518     1626-1637.

519     18.    Meisel, S., et al., Raman spectroscopy as a potential tool for detection of Brucella spp. in milk.

520     Appl. Environ. Microbiol., 2012. 78(16): p. 5575-5583.

521     19.    Meisel, S., et al., Identification of meat-associated pathogens via Raman microspectroscopy. Food

522     microbiology, 2014. 38: p. 36-43.

523     20.    Majed, N., et al., Identification of functionally relevant populations in enhanced biological

524     phosphorus removal processes based on intracellular polymers profiles and insights into the metabolic

525     diversity and heterogeneity. Environmental science technology, 2012. 46(9): p. 5010-5017.

526     21.    Kusić, D., et al., Identification of water pathogens by Raman microspectroscopy. Water research,

527     2014. 48: p. 179-189.

528     22.    Xu, J., et al., Label-free discrimination of Rhizobial bacteroids and mutants by single-cell Raman

529     microspectroscopy. Analytical chemistry, 2017. 89(12): p. 6336-6340.

530     23.    Li, Y., et al., Toward Better Understanding of EBPR Systems via Linking Raman-Based

531     Phenotypic Profiling with Phylogenetic Diversity. Environmental science technology, 2018. 52(15): p.

532     8596-8606.

533     24.    Große, C., et al., Label-free imaging and spectroscopic analysis of intracellular bacterial

534     infections. Analytical chemistry, 2015. 87(4): p. 2137-2142.

535     25.    Majed, N., et al., Evaluation of intracellular polyphosphate dynamics in enhanced biological

536     phosphorus removal process using Raman microscopy. Environmental science technology, 2009. 43(14):

537     p. 5436-5442.

538     26.    Ji, Y., et al., Raman spectroscopy provides a rapid, non‐invasive method for quantitation of starch

539     in live, unicellular microalgae. Biotechnology journal, 2014. 9(12): p. 1512-1518.

540    27.    Wang, D., et al., Side-stream enhanced biological phosphorus removal (S2EBPR) process

541    improves system performance-A full-scale comparative study. Water research, 2019. 167: p. 115109.

542    28.    He, Y., et al., Label-free, simultaneous quantification of starch, protein and triacylglycerol in single

543    microalgal cells. Biotechnology for biofuels, 2017. 10(1): p. 275.

544    29.    Xie, C., et al., Identification of single bacterial cells in aqueous solution using confocal laser

545    tweezers Raman spectroscopy. Analytical chemistry, 2005. 77(14): p. 4390-4397.

546    30.    Stöckel, S., et al., Raman spectroscopic detection and identification of Burkholderia mallei and

547    Burkholderia pseudomallei in feedstuff. Analytical bioanalytical chemistry, 2015. 407(3): p. 787-794.

548    31.    Ali, N., et al., Sample-size planning for multivariate data: a Raman-spectroscopy-based example.

549    Analytical chemistry, 2018. 90(21): p. 12485-12492.

550    32.    Beleites, C., et al., Sample size planning for classification models. Analytica chimica acta, 2013.

551    760: p. 25-33.

552    33.    Rényi, A. On measures of entropy and information. in Proceedings of the Fourth Berkeley

553    Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of

554    Statistics. 1961. The Regents of the University of California.

555    34.    Müller, A., Integral probability metrics and their generating classes of functions. Advances in

556    Applied Probability, 1997: p. 429-443.

557    35.    Kullback, S. and R.A. Leibler, On information and sufficiency. The annals of mathematical

558    statistics, 1951. 22(1): p. 79-86.

559    36.    Gretton, A., et al. Measuring Statistical Dependence with Hilbert-Schmidt Norms. 2005. Berlin,

560    Heidelberg: Springer Berlin Heidelberg.

561    37.    Jolliffe, I.T., Principal Components in Regression Analysis, in Principal Component Analysis.

562    1986, Springer New York: New York, NY. p. 129-155.

563    38.    Schölkopf, B., A. Smola, and K.-R. Müller, Nonlinear Component Analysis as a Kernel

564    Eigenvalue Problem. Neural Computation, 1998. 10(5): p. 1299-1319.

565    39.    Gu, A.Z., et al., Functionally relevant microorganisms to enhanced biological phosphorus removal

566    performance at full‑scale wastewater treatment plants in the United States. Water Environment

567    Research, 2008. 80(8): p. 688-698.

568    40.    Majed, N. and A.Z. Gu, Application of Raman microscopy for simultaneous and quantitative

569    evaluation of multiple intracellular polymers dynamics functionally relevant to enhanced biological

570    phosphorus removal processes. Environmental science technology, 2010. 44(22): p. 8601-8608.

571    41.    Onnis‑Hayden, A., et al., Impact of solid residence time (SRT) on functionally relevant microbial

572    populations and performance in full‑scale enhanced biological phosphorus removal (EBPR) systems.

573    Water Environment Research, 2019. 92(3): p. 389-402.

574    42.    Onnis‑Hayden, A., et al., Survey of full‑scale sidestream enhanced biological phosphorus

575    removal (S2EBPR) systems and comparison with conventional EBPRs in North America: Process

576    stability, kinetics, and microbial populations. Water Environment Research, 2019. 92(3): p. 403-417.

577    43.    De Gelder, J., et al., Reference database of Raman spectra of biological molecules. Journal of

578    Raman Spectroscopy, 2007. 38(9): p. 1133-1147.

579    44.    Chang, H. and D.-Y. Yeung, Robust path-based spectral clustering. Pattern Recognition, 2008.

580    41(1): p. 191-203.

581

582

583

584

585

586