

1 Two-target quantitative PCR to predict library composition for shallow shotgun  
2 sequencing

3

4 Matthew Y. Cho<sup>1,2</sup>, Marc Oliva<sup>3,4</sup>, Anna Spreafico<sup>3</sup>, Bo Chen<sup>5</sup>, Xu Wei<sup>5</sup>, Yoojin Choi<sup>1,2</sup>,  
5 Rupert Kaul<sup>1,2</sup>, Lillian L. Siu<sup>3</sup>, Bryan Coburn<sup>1,2\*</sup> and Pierre H. H. Schneeberger<sup>1,2\*</sup>

6

7 \* Co-senior Authors

8

9 1 Departments of Medicine and Laboratory Medicine & Pathobiology, University of  
10 Toronto, Toronto, Canada,

11 2 Department of Medicine, Division of Infectious Diseases, University Health Network,  
12 Toronto, Canada

13 3 Division of Medical Oncology and Hematology, Princess Margaret Cancer Centre,  
14 University of Toronto, Canada

15 4 Department of Medical Oncology, Catalan Institute of Oncology (Hospital Duran i  
16 Reynals), IDIBELL, Barcelona (Spain)

17 5 Department of Biostatistics, Princess Margaret Cancer Centre, University of Toronto,  
18 Canada

19 § Corresponding authors: Department of Medicine, Division of Infectious Diseases,  
20 University Health Network, M5G 1L7, Toronto, Canada.

21 Tel.: +41 61 284 8323, E-Mail: pierre.schneeberger@swisstph.ch Tel.: +1 416 581-7457,  
22 E-Mail: bryan.coburn@utoronto.ca

23

24 E-mail addresses:

25 MC matthewym.cho@mail.utoronto.ca

26 MO Moliva@iconcologia.net

27 AS anna.spreafico@uhn.ca

28 BC bo.chen@uhnresearch.ca

29 XW wei.xu@uhnresearch.ca

30 YC yoojinc.choi@mail.utoronto.ca

31 RK rupert.kaul@utoronto.ca

32 LLS Lillian.Siu@uhn.ca

33 BC bryan.coburn@utoronto.ca

34 PHHS pierre.schneeberger@swisstph.ch

35

36 Keywords:

37 Shotgun sequencing; shallow shotgun; microbiome; sample composition; host DNA  
38 proportion; metagenomics.

39 **Abstract**

40           Shotgun sequencing enables retrieving high resolution information from complex  
41    microbial communities. However, the technique is limited by missing information about  
42    host-to-microbe ratios observed in different sample types. This makes it challenging to  
43    plan sequencing experiments, especially in the context of high sample multiplexing  
44    and/or limited sequencing output. We evaluated a qPCR-based assay to predict host-to-  
45    microbe ratio prior to sequencing. Using a two-target assay aimed at conserved human  
46    and bacterial genes, we predicted human-to-microbe ratios in two sample types and  
47    validated it on independently collected samples. The assay enabled accurate prediction  
48    for a broad range of sample compositions.

49           **Introduction**

50           Shotgun sequencing allows interrogation of the metagenomic composition of  
51    ecological niches and has been increasingly utilized to characterize human-associated  
52    microbial communities. Shallow shotgun sequencing – sequencing to a per-sample read  
53    depth of  $10^5$  to  $10^6$  reads – provides taxonomic resolution greater than 16S amplicon  
54    sequencing and functional characterization of metagenomes, while being less  
55    expensive than whole genome sequencing or deep sequencing (typically  $10^7$  to  $10^9$   
56    reads/sample) (1). However, there is a trade-off between cost and adequacy which is  
57    especially problematic for samples of variable ratios of host to microbial DNA, where  
58    microbial reads may be displaced by human reads in a mixed sample (2). While this is  
59    generally not a concern for samples with high bacterial load, such as stool samples,  
60    samples with low or variable microbial DNA relative to human DNA are common in other  
61    regions of the body, such as the lung, nasopharynx, stomach, and duodenum (2, 3, 4,  
62    5). Microbial taxonomic and functional analyses of metagenomic data require sufficient

63 reads to draw robust conclusions. The ability to predict the proportion of microbial reads  
64 prior to sequencing would allow researchers to customize sequencing strategies for  
65 desired analyses, while optimizing the cost and time spent on metagenomic sequencing.

66 In this study, we used quantitative PCR to predict the ratio of human to microbial  
67 reads obtained from sequencing using three targets: the 16S rRNA gene, 18S rRNA  
68 gene, and human beta-actin (ACTB) to quantitate DNA of bacterial, fungal, or human  
69 origin, respectively (6-8). We compared the ratios of bacterial to human DNA  
70 determined via qPCR to the percent microbial/human DNA determined via shallow  
71 shotgun sequencing in samples with variable bacterial DNA. We derived a prediction  
72 model from oropharyngeal swabs and stool samples, and evaluated it in a set of  
73 independently collected samples, including rectal swabs and vaginal secretion samples.  
74 Finally, we generated an easy-to-use tool based on qPCR data to predict sample  
75 composition and sequencing depth required given a desired analytical outcome.

## 76 **Results and Discussion**

77 To assess the impact of shallowing sequencing depth on different bacterial DNA  
78 proportions, we rarefied shotgun sequencing data from 4 sample types – stool,  
79 oropharyngeal, rectal, and vaginal – to depths of 1000 to 1 million reads/sample. We  
80 then determined the alpha diversity of each rarefaction using three metrics: richness,  
81 Shannon index, and Berger-Parker index. Alpha diversity decreased in a sample type-  
82 specific manner as sequencing depth decreased (**Fig. 1**). Notably, while vaginal  
83 samples have the lowest alpha diversity in all three metrics of the four sample types,  
84 alpha diversity decreased at the slowest rate as sequencing depth decreased (**Fig. 1**).

85 Conversely, while rectal swab samples had similar Shannon index and Berger-Parker  
86 index values at  $10^6$  microbial reads to oropharyngeal and stool samples, alpha diversity  
87 in rectal samples diminished at a greater rate as sequencing depth decreased (**Fig. 1B-C**). Since this effect is sample type-specific, it is critical to predict sample composition *a*  
88 *priori* to ensure sufficient reads for the desired analysis for the given sample type.

90 qPCR is a widespread and robust technique available in many molecular biology  
91 laboratories. Its availability as well as cheap associated costs, especially compared to  
92 experiments involving high-throughput sequencing techniques, makes it an ideal  
93 candidate to use to predict sample composition prior to sequencing. In this study, we  
94 assessed the potential of qPCR to predict sample-specific ratios of human to microbe  
95 DNA using different amplification targets. Using a multivariate approach, 5 models were  
96 generated mapping 16S rRNA gene, 18S rRNA gene, and human beta-actin (ACTB)  
97 qPCR-derived cycle thresholds (Ct) to observed percentage of microbial reads for a  
98 sample set consisting of oropharyngeal swabs and stool samples. Microbial reads were  
99 defined as any read which did not align/match with a human genome reference. The  
100 following models were tested: (A) a linear fit using 16S rRNA gene and ACTB Ct values,  
101 (B) a linear fit using 16S rRNA gene, 18S rRNA gene, and ACTB Ct values, (C) a linear  
102 fit using logit transformed 16S rRNA gene and ACTB Ct values, (D) a linear fit using  
103 logit transformed 16S rRNA gene, 18S rRNA gene, and ACTB Ct values, and (E) a  
104 nonlinear regression model based on the logistic growth equation using 16S rRNA gene  
105 and ACTB Ct values (**Supplementary figure 1A**). We compared goodness-of-fit for  
106 each model and observed  $R^2$ -values of 0.880, 0.880, 0.920, 0.920, and 0.990 for  
107 models A – E, respectively (**Figure 2A, Supplementary figure 1A**). Observed residuals

108 had a min-max range of 67.56, 68.50, 58.93, 59.07, and 42.61 for models A – E,  
109 respectively (**Supplementary figure A**). Based on these findings, model E turned out to  
110 be the best fitting model to predict sample composition using qPCR, with an equation  
111 of % microbial reads =  $(2.7201549)/((99.50267)*e^{-0.7218*(ACTB-16S)})+ 0.02733$ . In  
112 addition, 18S rRNA Ct value was not found to be an informative predictor and was  
113 hence removed from the model. In **Figure 2B**, we show the goodness-of-fit and  
114 residuals observed with model E across the range of qPCR differences (-8.16% to  
115 +34.45%). We observed homogeneous fit and variance indicating that the model  
116 performs well for all observed host to microbe DNA ratios. However, we also observed  
117 that the model loses accuracy at each end of the range due to the s initial dataset used  
118 and sigmoidal curve generated, with limits approximately at 4% and 98%. This bias is  
119 likely introduced at different steps of the process. For instance, sequencing error, and  
120 resulting false negative and positive hits when mapping reads to the human database  
121 are likely to account for this bias. Another potential source of bias could be introduced  
122 by the carryover of contaminants between sequencing runs, hence resulting in a  
123 composition change which is not picked up by the qPCR conducted *a priori*.

124 Using the equation derived from model E, we evaluated our approach on two different,  
125 independently collected sample types including vaginal secretions and rectal swabs. In  
126 **Fig 2C**, we show the relation between observed microbial reads percentages and the  
127 difference in Ct between 16S and ACTB qPCR, derived from our validation dataset,  
128 alongside a curve of expected values derived from model E. We observed the  
129 difference between predicted and observed microbial reads percentages to range from -  
130 18.80% to +19.22% with a mean of +0.944% (**Supplementary figure 1B**). In **Fig 2D**,

131 we show that this difference is consistent across the range of observed % microbial  
132 reads. Compared to the other models, model E best described the validation dataset,  
133 with a median difference of 0.25% and a standard deviation of 9.10% (**Supplementary**  
134 **table 1B**). For comparison, model E described the initial sample set of  
135 oropharyngeal/stool samples with a median difference of 0.14% and a standard  
136 deviation of 4.35% (**Supplementary table 1A**). Since the model performed similarly  
137 between the two datasets, we concluded that the model was able to describe a relation  
138 between 16S and  $\beta$ -actin qPCR and shotgun sequencing metagenomic data in a  
139 sample type-independent manner for microbial densities between 4% and 98%. We then  
140 developed a tool based on our model and the rarefaction curves on different samples  
141 type which predicts % microbial reads based on qPCR data and suggests a target  
142 number of reads based on sample type and desired analysis (**Supplementary**).  
  
143 The limitations of our study are as follows: The samples used in our study were low in  
144 fungal content. Therefore, our model may not accurately predict microbial content in  
145 sample sets where the majority of samples are rich in fungal content.  
  
146 Moreover, as our results are based on protocols using specific reagents and  
147 technologies for both sequencing and qPCR, our tool may not accurately predict  
148 sequencing results when protocols, reagents, and/or technologies differ. However,  
149 given that we have established a robust link among 16S qPCR,  $\beta$ -actin qPCR, and  
150 sample content by sequencing, our approach can be easily adapted to fit different  
151 experimental settings.

152 **Conclusion**

153 We have shown that shallowing shotgun sequencing depth can reduce measured alpha  
154 diversity in all measured sample types, with more diverse communities being more  
155 strongly negatively affected. We found that qPCR can function as a predictive tool for  
156 sample composition that was strongly correlated with shotgun sequencing data. We  
157 were able to create a model that can describe and predict variable sample types. We  
158 hope that our tool and methodology may help fellow researchers screen for  
159 sequenceable samples or allow for better optimization of sequencing.

160 **Methods**

161 *qPCR*

162 Samples were probed separately for the 16S rRNA gene, the 18S rRNA gene, and the  
163 human  $\beta$ -actin gene. All reactions were conducted in duplicate and RNase-free water  
164 was used as negative control. Each well contained 2  $\mu$ L of sample DNA, 5  $\mu$ L of  
165 Taqman Universal PCR mix (Applied Biosystems, Foster City, CA), 0.3  $\mu$ M of forward  
166 primer, 0.3  $\mu$ M of reverse primer, and 0.2  $\mu$ M of primer probe. PCR was performed on a  
167 QuantStudio 6 Flex (Thermo Fisher Scientific, Waltham, MA) platform. Cycling was  
168 done as follows: 10 minutes at 95°C followed by 45 cycles of 95°C for 15 seconds and  
169 60°C for 1 minute.

170 For 16S qPCR, we used forward primer “TCCTACGGGAGGCAGCAGT” (Invitrogen,  
171 Carlsbad, CA) and reverse primer “GGACTACCAGGGTATCTAACCTGTT” (Invitrogen,  
172 Carlsbad, CA).(3) We used a FAM probe “CGTATTACCGCGGCTGCTGGCAC” with  
173 NFQ-MGB quencher (Applied Biosystems, Foster City, CA).(3)

174 For 18S qPCR, we used forward primer “GGRAAACTCACCAAGGTCCAG” (Integrated  
175 DNA Technologies, Coralville, IA) and reverse primer “GSWCTATCCCCAKCACGA”  
176 (Integrated DNA Technologies, Coralville, IA).(1) We used a FAM probe  
177 “TGGTGCATGGCCGTT” with NFQ-MGB quencher (Applied Biosystems, Foster City,  
178 CA).(7)

179 For human qPCR, we used a  $\beta$ -actin gene specific forward primer  
180 “CGGCCTTGGAGTGTATTAAAGTA” (Invitrogen, Carlsbad, CA) and reverse primer  
181 “TGCAAAGAACACGGCTAAGTGT” (Invitrogen, Carlsbad, CA).(5) We used a VIC  
182 probe “TCTGAACAGACTCCCCATCCCAAGACC” with 3QSY quencher (Applied  
183 Biosystems, Foster City, CA).(8)

184 *Library preparation and sequencing*

185 Libraries were prepared using Nextera Flex (Illumina, San Diego, CA) kits with the  
186 Nextera XT indices (Illumina, San Diego, CA). Barcoded sample libraries were pooled  
187 together to a concentration of 17.6 ng/ul which measured with a high-sensitivity DNA  
188 assay on a Qubit (Thermo Fisher Scientific, Waltham, MA) platform. A Mid-output  
189 reagent kit (Illumina, San Diego, CA) was used to sequence on the Miniseq, while a SP  
190 reagent kit (Illumina, San Diego, CA) was used on the Novaseq platform, both in  
191 2x150bp mode.

192 *Read filtering and Taxonomic profiling*

193 We filtered human reads from non-human reads using KneadData based on a human  
194 genome index for Bowtie 2 (9, 10). We considered sequence reads that did not match  
195 the database as microbial reads in our analyses. Taxonomic annotation was conducted

196 using MetaPhlAn 2.0 and the ChocoPhlAn database (11). Rarefactions were performed  
197 using seqtk-1.3 to subsample the microbial reads of individual samples (12). Subsample  
198 compositions will be identified using MetaPhlAn2, and OTU tables were generated (11).  
199 Diversity indexes were calculated using Past 4 (13).

200 *Model generation*

201 We used XLSTAT version 2019.4.2 (Addinsoft Inc., New York, NY) to generate  
202 multivariate linear regressions using either 16S and ACTB qPCR cycle and microbial  
203 reads percentages (Models A and C) or 16S, 18S, and ACTB qPCR cycle thresholds  
204 and microbial reads percentages (Models B and D). Multivariate linear regressions  
205 (models C and D) were also performed following a logit transformation of microbial  
206 reads percentages. Finally, for model E, we generated the non-linear regression model  
207 using the logistic growth equation in GraphPad Prism version 8.3.0 for Windows  
208 (GraphPad Software, San Diego, CA).

209 **Figures**

210 **Figure 1. Alpha diversity indices are shown across a range of simulated**  
211 **sequencing depths from 1E3 to 1E6 reads per sample.** (A) Sample-specific  
212 rarefaction curves of species richness. (B) Shannon index calculated across a range of  
213 rarefactions, by sample type. (C) Sample dominance, measured with the Berger-Parker  
214 index, across a range of sequencing depths, stratified by sample type.

215 **Figure 2. Statistical model to predict sample composition using qPCR prior to**  
216 **high-throughput sequencing (A)** Sigmoidal model generated from oropharyngeal  
217 swabs and stool samples depicting the relationship between the difference of human

218 (ACTB) and bacterial (16S) qPCR values (Ct) with the percentage of microbial reads ( $R^2$   
219 =0.990). Nonlinear regression line (solid) is based on the following logistic growth  
220 equation: % microbial reads = (2.7201549)/((99.50267)\*e^(-0.7218\*(ACTB-16S))+  
221 0.02733). One-tailed 95% prediction interval is depicted with a dotted line. **(B)** Model  
222 residuals. **(C)** Fitting of validation sample set on prediction model. The orange dots  
223 represent values derived from a validation sample set composed of vaginal secretions  
224 and rectal swabs samples and correlate well ( $R^2$  = 0.930) with the prediction model  
225 (solid black line). **(D)** Difference between expected and observed composition across  
226 the range of microbial content.

227 **Supplementary Figure 1. (1A)** Residuals for 5 multivariate models generated using a  
228 sample set comprised of oropharyngeal swabs and stool samples. i) Model A  
229 represents a linear fit taking into account microbial and human-derived qPCR values; ii)  
230 model B represents a linear fit taking into account microbial, fungal, and human-derived  
231 qPCR values; iii) model C represents a linear fit taking into account microbial and  
232 human-derived qPCR values after a logit transform of the data; iv) model D represents a  
233 linear fit taking into account microbial, fungal, and human-derived qPCR values after a  
234 logit transform of the data; and v) model E represents a nonlinear regression model  
235 based on the logistic growth equation taking into account microbial and human-derived  
236 qPCR values. Error bars depict 1 standard deviation centered around the mean. **(1B)**  
237 Difference between observed and predicted percentage of microbial reads, by model,  
238 using a validation dataset comprised of independently collected rectal swabs and  
239 vaginal secretion samples.

240 **Supplementary Table 1. (1A)** Table summarizing residual values and model type of  
241 the 5 statistical models (A-E) tested in this study. **(1B)** Residual values of the 5 models  
242 when applied to an independent, validation dataset.

243 **References**

244 1. Hillmann, B., G. A. Al-ghalith, R. R. Shields-cutler, Q. Zhu, D. M. Gohl, K. B. Beckman, R. Knight, and D.  
245 Knights. 2018. Evaluating the Information Content of Shallow Shotgun Metagenomics. *mSystems* 3: 1–12.

246 2. Leo, S., N. Gaïa, E. Ruppé, S. Emonet, M. Girard, V. Lazarevic, and J. Schrenzel. 2017. Detection of  
247 Bacterial Pathogens from Broncho-Alveolar Lavage by Next-Generation Sequencing. *International*  
248 *journal of molecular sciences* 18: 1–13.

249 3. Nadkarni, M. A., F. E. Martin, N. A. Jacques, and N. Hunter. 2002. Determination of bacterial load by  
250 real-time PCR using a broad-range (universal) probe and primers set. *Microbiology* 148: 257–266.

251 4. Biesbroek G., E. A. Sanders, G. Roeselers, M. P. M. Caspers, K. Trzciński, D. Bogaert, and B. J. F. Keijser.  
252 2012. Deep sequencing analyses of low density microbial communities: working at the boundary of  
253 accurate microbiota detection. *PLoS One* 7(3)

254 5. Bogaert, D., B Keijser, S Huse, J. Rossen, R. Veenhoven, E. van Gils, J. Bruin, R. Montijn, M. Bonten, and  
255 E. Sanders. 2011. Variability and Diversity of Nasopharyngeal Microbiota in Children: A Metagenomic  
256 Analysis. *PLoS Biol* 6(2)

257 6. Sender, R., S. Fuchs, and R. Milo. 2016. Revised Estimates for the Number of Human and Bacteria Cells  
258 in the Body. *PLoS Biol* 14(8)

259 7. Liu, C. M., S. Kachur, M. G. Dwan, A. G. Abraham, M. Aziz, P. R. Hsueh, Y. T. Huang, J. D. Busch, L. J.  
260 Lamit, C. A. Gehring, P. Keim, and L. B. Price. 2012. FungiQuant: a broad-coverage fungal quantitative  
261 real-time PCR assay. *BMC microbiology* .

262 8. Hasan, M. R., A. Rawat, P. Tang, P. v. Jithesh, E. Thomas, R. Tan, and P. Tilley. 2016. Depletion of  
263 human DNA in spiked clinical specimens for improvement of sensitivity of pathogen detection by next-  
264 generation sequencing. *Journal of Clinical Microbiology* 54: 919–927.

265 9. Langmead, B., and S. L. Salzberg. 2012. Fast gapped-read alignment with Bowtie 2. *Nature Methods* 9:  
266 357–359.

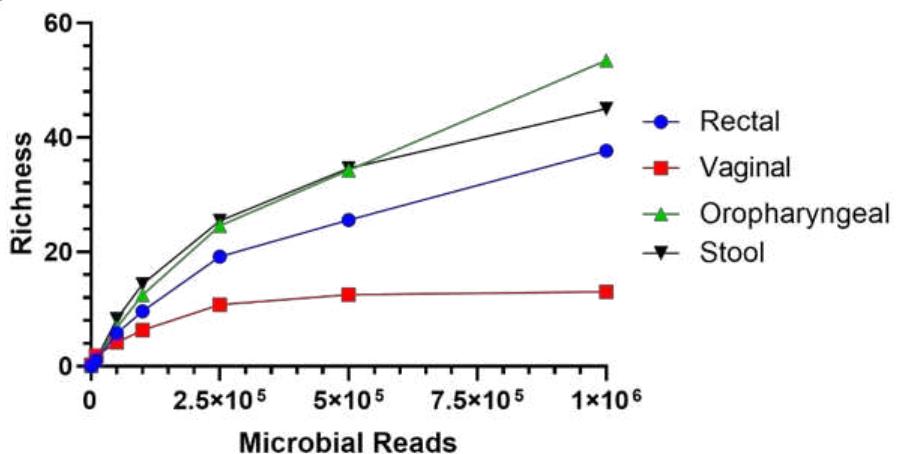
267 10. Huttenhower, C. KneadData | The Huttenhower Lab. .

268 11. Segata, N., L. Waldron, A. Ballarini, V. Narasimhan, O. Jousson, and C. Huttenhower. 2013.  
269 MetaPhiAn -1-. *Nat Methods* 9: 811–814.

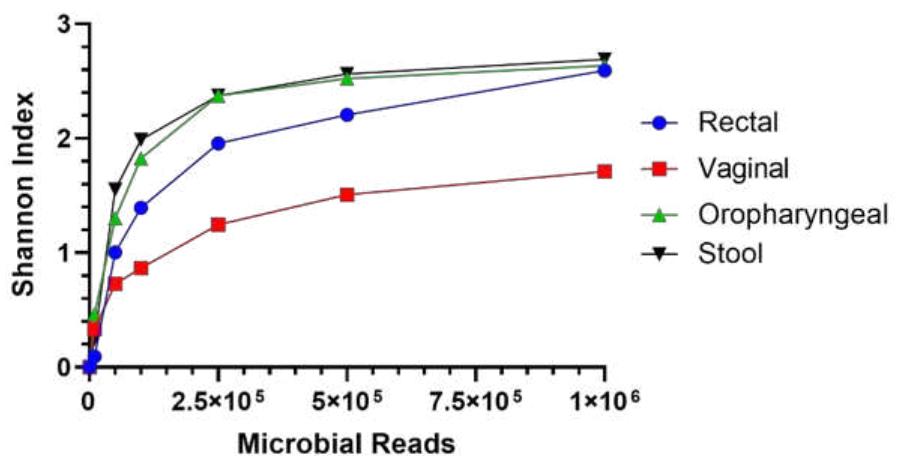
270 12. Li, H. GitHub - lh3/seqtk: Toolkit for processing sequences in FASTA/Q formats.

271 13. Hammer, Ø., Harper, D.A.T., Ryan, P.D. 2001. PAST: Paleontological statistics software package for  
272 education and data analysis. *Palaeontologia Electronica* 4(1): 9pp.

(A)



(B)



(C)

