

FOCUS: an Alignment-free Model to Identify Organisms in Metagenomes Using Non-negative Least Squares

One of the major goals in metagenomics is to identify the organisms present in a microbial community from unannotated shotgun sequencing reads. Taxonomic profiling has valuable applications in biological and medical research, including disease diagnostics. Most currently available approaches do not scale well with increasing data volumes, which is important because both the number and lengths of the reads provided by sequencing platforms keep increasing. Here we introduce FOCUS, an agile composition based approach using non-negative least squares (NNLS) to report the focal organisms present in metagenomic samples and profile their abundances. FOCUS was tested with simulated and real metagenomes, and the results show that our approach accurately predicts the organisms present in microbial communities. FOCUS was implemented in Python. The source code and web-server are freely available at <http://edwards.sdsu.edu/FOCUS>.

1
2 **FOCUS: an Alignment-free Model to Identify Organisms in Metagenomes Using Non-**
3 **negative Least Squares**

4

5 Genivaldo Gueiros Z. Silva¹, Daniel A. Cuevas¹, Bas E. Dutilh ^{4, 5}, and Robert A. Edwards ^{1, 2, 3, 5, 6}

6 *

7 ¹Computational Science Research Center, ²Department of Computer Science, and ³Department of
8 Biology, San Diego State University, 5500 Campanile Drive, San Diego, CA 92182, USA,

9 ⁴Centre for Molecular and Biomolecular Informatics, Radboud Institute for Molecular Life
10 Sciences, Radboud University Medical Centre, Geert Grooteplein 28, 6525 GA, Nijmegen, The
11 Netherlands, ⁵Department of Marine Biology, Institute of Biology, Federal University of Rio de
12 Janeiro, Brazil, ⁶Division of Mathematics and Computer Science, Argonne National Laboratory,
13 9700 S. Cass Ave, Argonne, IL 60439, USA

14

15 *For correspondence please contact Dr. Robert Edwards at redwards@mail.sdsu.edu.

16

17 **Abstract**

18 One of the major goals in metagenomics is to identify the organisms present in a microbial
19 community from unannotated shotgun sequencing reads. Taxonomic profiling has valuable
20 applications in biological and medical research, including disease diagnostics. Most currently
21 available approaches do not scale well with increasing data volumes, which is important because
22 both the number and lengths of the reads provided by sequencing platforms keep increasing. Here we
23 introduce FOCUS, an agile composition-based approach using non-negative least squares (NNLS) to
24 report the focal organisms present in metagenomic samples and profile their abundances. FOCUS
25 was tested with simulated and real metagenomes, and the results show that our approach accurately
26 predicts the organisms present in microbial communities. FOCUS was implemented in Python. The
27 source code and web-server are freely available at <http://edwards.sdsu.edu/FOCUS>.

28

29 **Introduction**

30 Microbes are more abundant than any other cellular organism (Whitman et al. 1998), and
31 it is important to understand which organisms are present and what they are doing (Handelsman
32 2004). In many environments a majority of the microbial community members cannot be
33 cultured, and metagenomics is a powerful tool to directly probe uncultured genomes and
34 understand the diversity of microbial communities by using only their DNA (Sharon and Banfield
35 2013).

36 Understanding microbial communities is important in many areas of biology. For
37 example, metagenomes can distinguish taxonomic and functional signatures of microbes
38 associated with marine animals (Trindade-Silva et al. 2012) or disease states (Belda-Ferre et al.
39 2012). Large sequencing volumes, short read lengths, and sequencing errors make the task of

40 identifying the diversity of organisms present in metagenomes challenging (Mande et al. 2012).
41 Many programs exist for this, and they are either homology- or composition-based.

42 Homology-based programs normally use BLAST program (Altschul et al. 1997) to
43 identify the best hit in a large database output. In MG-RAST (Meyer et al. 2008) sequences are
44 aligned to a set of databases in order to classify the metagenomic sample. MetaPhlAn (Segata et
45 al. 2012) and GenomePeek (McNair and Edwards) use a reduced database containing only
46 marker genes, e.g., unique clades and housekeeping genes, allowing the BLAST search to be fast.
47 PhymmBL (Brady and Salzberg 2011) improves the BLAST results using interpolated Markov
48 models. GASiC (Lindner and Renard 2013) uses Bowtie (Langmead et al. 2009) and the
49 reference genomes similarities to correct the observed abundance estimated. Parallel-Meta (Su et
50 al. 2012) a fast program, which requires a GPU, uses megaBLAST (Zhang et al. 2000) and HMM
51 (Hidden Markov Model) to improve the homology result. Most of these applications classify
52 sequences individually, and generate a taxonomic profile by summing the bins.

53 In general, composition-based approaches use oligonucleotide (k -mer) frequencies. Taxy
54 (Meinicke et al. 2011) uses oligonucleotide distribution in metagenomes and in reference
55 genomes and uses mixture modeling to identify the organisms present in the metagenome, and
56 RAiphy (Nalbantoglu et al. 2011) identifies organisms using oligonucleotides and relative
57 abundance index.

58 We developed a new approach that reconstructs a taxonomic profile using an ensemble k -
59 mer composition of the entire metagenome. We compute the optimal set of organism abundances
60 using non-negative least squares (NNLS) to match the metagenome k -mer composition to
61 organisms in a reference database. K -mers have previously been used to cluster unknown
62 sequences (Teeling et al. 2004; McHardy et al. 2007) and NNLS has been used to identify the
63 genera present in metagenomic samples based on variations in gene count (Carr et al. 2013). Here
64 we combine these two approaches in FOCUS, an ultra fast, accurate, composition based approach

65 to identify the taxa present in a metagenome. We compare the performance of FOCUS to GASiC,
66 MetaPhlAn, RAphy, PhymmBL, Taxy, and MG-RAST.

67

68 Methods

69 FOCUS workflow is described in Figure 1. As in most composition-based approaches, a
70 training set is pre-generated using the complete genomes information, and here the non-negative
71 least squares (NNLS) is applied to compute the relative abundance of each organism in the
72 database into the unknown data.

73

74 Reference dataset

75 FOCUS requires a group of reference genomes to model and identify the organisms
76 present in a metagenome. 2,766 complete genomes were downloaded from the SEED servers
77 (Aziz et al. 2012) on 20 December 2013 (see Supplementary Table 1). k -mer frequencies ($k=6-8$,
78 default: $k=7$) were calculated for both strands using Jellyfish 1.1.6 (Marçais and Kingsford 2011),
79 reducing the number of dimensions (Strous et al. 2012), and k -mer counts were normalized by the
80 sum of frequencies. The user can also create their own training set, which is scalable to the
81 quickly increasing number of available reference genomes because it also uses Jellyfish in the k -
82 mer counting.

83

84 Simulated and real metagenomes

85 In order to evaluate FOCUS performance, a simulated dataset of short sequences
86 (SimShort), containing 500,000 single 100 nt reads, was created using the supplied error model
87 for Illumina GA IIx with TrueSeq SBS Kit v5–GA using GemSim (McElroy et al. 2012)
88 (Supplementary Table 2). The previous published high complexity simulated dataset (SimHC)
89 from FAMeS (Mavromatis et al. 2007) was also used in the evaluation. Moreover, real human

90 oral sample datasets were selected as test cases: one under healthy conditions and one under
91 disease conditions (MG-RAST accession 4447943.3 and 4447192.3) (Belda-Ferre et al. 2012)
92 were selected as a test case.

93

94 Non-Negative Least Squares (NNLS)

95 The estimation of a parameterized model to understand some data is a fundamental
96 problem in data modeling. Nevertheless, the estimation is not always easy, e.g., in problems like
97 metagenome profiling that cannot have negative values for the fitted parameters. In such case, a
98 solution can be estimated using NNLS, which is defined as:

99 Given a matrix $A \in \mathbb{R}^{m \times n}$ and a vector $b \in \mathbb{R}^m$, where $m \geq n$, find a non-negative vector $x \in$
100 \mathbb{R}^n to minimize the function (1).

101

102

(1)

$$103 f(x) = \frac{1}{2} \|Ax - b\|^2, \text{ where } x \geq 0 \text{ and } \sum_{i=1}^n x_i = 1$$

104 In FOCUS, the reference matrix A is composed of m k -mer frequencies from n genomes,
105 while a vector describing the user's metagenomic dataset is calculated from the k -mer frequencies
106 of both strands from the dataset using Jellyfish. FOCUS uses non-negative least squares to
107 compute the set of k -mer frequencies x that explains the optimal possible abundance of k -mers in
108 the user's metagenome by selecting the optimal number of frequencies from the matrix A . We
109 minimize the sum of squared differences (1) using the open source Scipy library (Jones et al.
110 2001) which has a module for the NNLS algorithm which solves the KKT (Karush-Kuhn-Tucker)
111 conditions (Lawson and Hanson 1987). We added Tikhonov regularization (Garda and Galias
112 2012) to deal with genomes that have similar k -mer compositions.

113

114 Jackknife resampling of the data

115 We implemented a jackknife resampling strategy to assess the robustness of the results.
116 50% of the reads were randomly resampled 1000x, and the species frequencies recalculated. For
117 each species, these 1000 frequencies were averaged and the standard deviation calculated to
118 estimate the spread.

119 **Web-based and graphical user interface version**

120 As an alternative to the command line version of the program, we have created a user-
121 friendly web version and a graphical user interface (GUI) for Microsoft Windows users. The web
122 server and the GUI are available at <http://edwards.sdsu.edu/FOCUS>.

123 **Results and Discussion**

124

125 **Evaluation and comparison with other tools**

126 All tools were run using default parameters and their default reference database, either
127 online (MG-RAST) or using one core on a server with 24 processors x 6 cores Intel(R) Xeon(R)
128 CPU X5650 @ 2.67GHz and 189 GB RAM. We only compared GASiC to the SimHC dataset
129 which had the results previously published (Lindner and Renard 2013). We tried to run the tool;
130 however, it requires a large amount of storage during the process to save its output data.

131 For the real data, two metagenomic datasets were selected. First, the metagenomic sample
132 of the human oral cavity from diseased conditions was used. MetaPhlAn apparently over
133 predicted the genera *Veillonella* due to the short genome, and Taxy did not predict *Prevotella* hits
134 (see Figure 2) as described in (Belda-Ferre et al. 2012). FOCUS was able to profile the organisms
135 in only 41 seconds. Taxy took about 45 seconds, MetaPhlAn took about 3 minutes, RAiphy took
136 52 minutes, MG-RAST took 3 days, and PhymmBL took 1 week and 6 days. Using random
137 subsets for the oral metagenome, we tested the tools scalability and showed that FOCUS and
138 Taxy profile metagenomes in constant time (see Figure 3).

139 The oral metagenome from the healthy condition was used. MetaPhlAn possibly over
140 predicted the genera *Neisseria*, and Taxy was not able to predict *Rothia* hits (see Figure 4).
141 FOCUS profiled the metagenome in only 35 seconds. Taxy took about 41 seconds, MetaPhlAn
142 took about 2 minutes, RAphy took 48 minutes, MG-RAST took 3 days, and PhymmBL took 9
143 days.

144 For the simulated data, we removed species from the reference dataset that are present in
145 this dataset and tried to predict the genera present in the SimShort dataset. A major limitation of
146 many of the approaches discussed here is that the underlying databases cannot be changed. Only
147 FOCUS, RAphy, GASiC, and PhymmBL allow the end user to change their reference database.
148 FOCUS and PhymmBL best predicted the correct genera while RAphy could not correctly
149 predict their abundance (see Figure 5). FOCUS had the fastest performance (45 seconds) in
150 comparison; RAphy took about 2 hours, while PhymmBL took approximately 5 days. See
151 Supplementary Figure 1 to 4 for other taxonomy resolutions.

152 For the SimHC simulated metagenomes, the genera present in the dataset were deleted
153 from the training dataset, and we evaluated the class-level prediction. The tested tools correctly
154 predicted the classes, except that RAphy over predicted the top two classes (see Figure 6).
155 Again, FOCUS was the fastest tool (30 seconds) in comparison to RAphy, which took about 1
156 hour and 50 minutes, and PhymmBL, which took about 4 days. See Supplementary Figure 5 to 7
157 for other taxonomy levels.

158 Furthermore, for the SimHC dataset, we ran all the previously used tools and the GASiC
159 published results to evaluate the genera-level prediction. GASiC and PhymmBL had the best
160 predictions, and FOCUS failed in the prediction of 4 minor genera probably because many
161 organisms present in the SimHC dataset were not included in the FOCUS database (see Figure 7).
162 We did not compare the running time because we extracted the GASiC results from its paper;

163 however, in the original paper it took 2 days and needed at least 500 GB of storage to analyze the
164 SimHC simulated metagenome.

165 The very small standard deviations observed after jackknife re-sampling indicate the
166 robustness of our results.

167 These tests were performed on a server; however, FOCUS is also ultra fast on a simple
168 computer. For example, we profiled the real dataset in 1 minute and 45 seconds using an Intel(R)
169 Core(TM) i3 @2.53 GHz and 1GB RAM.

170

171 **Limitations**

172 As with other methods created to profile metagenome sequences, FOCUS depends on a
173 curated database of microbial reference genomes in order to predict a specific genus. If a
174 reference genome is absent, the tool will predict the closest reference available.

175

176 **Conclusions**

177 Here we present FOCUS, an agile solution to identify the organisms present in
178 metagenomic samples that does not rely on mapping individual reads, but instead determines the
179 taxonomic composition of the entire metagenome at once by using NNLS. This makes FOCUS
180 an extremely fast and scalable solution to profile the focal taxa in a metagenome. FOCUS reports
181 very similar species compositions as currently available, state of the art metagenome profiling
182 tools.

183

184 **Availability and requirements**

185 Project name: FOCUS

186 Project and web server home page: <http://edwards.sdsu.edu/FOCUS>

187 Operating system: the program has a command line version that works on OS X and Unix, and a

188 GUI for Microsoft Windows users.

189 Programming language: Python 2.7.

190 Other requirements: Numpy (<http://www.numpy.org>), Scipy (<http://scipy.org>), Jellyfish

191 (<http://www.cbcn.umd.edu/software/jellyfish>), and Python programming language

192 (<http://www.python.org>).

193 License: GNU GPL3.

194 Any restrictions to use by non-academics: no special restrictions.

195

196 **Acknowledgements:**

197 We thank Dr. Peter Blomgren for help with the Advanced Numerical Analysis.

198

199 **References**

- 200 Altschul, S. F., T. L. Madden, A. A. Schäffer, J. Zhang, Z. Zhang, W. Miller, and D. J. Lipman,
201 1997: Gapped BLAST and PSI-BLAST: a new generation of protein database search
202 programs. *Nucleic Acids Res.*, **25**, 3389–3402, doi:10.1093/nar/25.17.3389.
- 203 Aziz, R. K., and Coauthors, 2012: SEED Servers: High-Performance Access to the SEED
204 Genomes, Annotations, and Metabolic Models. *PLoS ONE*, **7**, e48053,
205 doi:10.1371/journal.pone.0048053.
- 206 Belda-Ferre, P., L. D. Alcaraz, R. Cabrera-Rubio, H. Romero, A. Simón-Soro, M. Pignatelli, and
207 A. Mira, 2012: The oral metagenome in health and disease. *ISME J.*, **6**, 46–56,
208 doi:10.1038/ismej.2011.85.
- 209 Brady, A., and S. Salzberg, 2011: PhymmBL expanded: confidence scores, custom databases,
210 parallelization and more. *Nat. Methods*, **8**, 367–367, doi:10.1038/nmeth0511-367.

- 211 Carr, R., S. S. Shen-Orr, and E. Borenstein, 2013: Reconstructing the Genomic Content of
212 Microbiome Taxa through Shotgun Metagenomic Deconvolution. *PLoS Comput Biol*, **9**,
213 e1003292, doi:10.1371/journal.pcbi.1003292.
- 214 Garda, B., and Z. Galias, 2012: Non-negative least squares and the Tikhonov regularization
215 methods for coil design problems. *2012 International Conference on Signals and*
216 *Electronic Systems (ICSES)*, 2012 International Conference on Signals and Electronic
217 Systems (ICSES), 1–5.
- 218 Handelsman, J., 2004: Metagenomics: Application of Genomics to Uncultured Microorganisms.
219 *Microbiol. Mol. Biol. Rev.*, **68**, 669–685, doi:10.1128/MMBR.68.4.669-685.2004.
- 220 Jones, E., T. Oliphant, and P. Peterson, 2001: SciPy: Open source scientific tools for Python.
221 <http://www.scipy.org/>, http://www.scipy.org/Citing_SciPy (Accessed October 23, 2013).
- 222 Langmead, B., C. Trapnell, M. Pop, and S. L. Salzberg, 2009: Ultrafast and memory-efficient
223 alignment of short DNA sequences to the human genome. *Genome Biol.*, **10**, R25,
224 doi:10.1186/gb-2009-10-3-r25.
- 225 Lawson, C., and R. J. Hanson, 1987: *Solving Least Squares Problems*. SIAM,.
- 226 Lindner, M. S., and B. Y. Renard, 2013: Metagenomic abundance estimation and diagnostic
227 testing on species level. *Nucleic Acids Res.*, **41**, e10, doi:10.1093/nar/gks803.
- 228 Mande, S. S., M. H. Mohammed, and T. S. Ghosh, 2012: Classification of metagenomic
229 sequences: methods and challenges. *Brief. Bioinform.*, **13**, 669–681,
230 doi:10.1093/bib/bbs054.
- 231 Marçais, G., and C. Kingsford, 2011: A fast, lock-free approach for efficient parallel counting of
232 occurrences of k-mers. *Bioinformatics*, **27**, 764–770, doi:10.1093/bioinformatics/btr011.
- 233 Mavromatis, K., and Coauthors, 2007: Use of simulated data sets to evaluate the fidelity of
234 metagenomic processing methods. *Nat. Methods*, **4**, 495–500, doi:10.1038/nmeth1043.
- 235 McElroy, K. E., F. Luciani, and T. Thomas, 2012: GemSIM: general, error-model based simulator
236 of next-generation sequencing data. *BMC Genomics*, **13**, 74, doi:10.1186/1471-2164-13-
237 74.
- 238 McHardy, A. C., H. G. Martín, A. Tsirigos, P. Hugenholtz, and I. Rigoutsos, 2007: Accurate
239 phylogenetic classification of variable-length DNA fragments. *Nat. Methods*, **4**, 63–72,
240 doi:10.1038/nmeth976.
- 241 McNair, K., and R. Edwards, GenomePeek – A tool for prokaryotic genome and metagenome
242 analysis.
- 243 Meinicke, P., K. P. Aßhauer, and T. Lingner, 2011: Mixture models for analysis of the taxonomic
244 composition of metagenomes. *Bioinformatics*, btr266, doi:10.1093/bioinformatics/btr266.
- 245 Meyer, F., and Coauthors, 2008: The metagenomics RAST server – a public resource for the
246 automatic phylogenetic and functional analysis of metagenomes. *BMC Bioinformatics*, **9**,
247 386, doi:10.1186/1471-2105-9-386.

- 248 Nalbantoglu, O. U., S. F. Way, S. H. Hinrichs, and K. Sayood, 2011: RAIPhy: Phylogenetic
249 classification of metagenomics samples using iterative refinement of relative abundance
250 index profiles. *BMC Bioinformatics*, **12**, 41, doi:10.1186/1471-2105-12-41.
- 251 Segata, N., L. Waldron, A. Ballarini, V. Narasimhan, O. Jousson, and C. Huttenhower, 2012:
252 Metagenomic microbial community profiling using unique clade-specific marker genes.
253 *Nat. Methods*, **9**, 811–814, doi:10.1038/nmeth.2066.
- 254 Sharon, I., and J. F. Banfield, 2013: Genomes from Metagenomics. *Science*, **342**, 1057–1058,
255 doi:10.1126/science.1247023.
- 256 Strous, M., B. Kraft, R. Bisdorf, and H. E. Tegetmeyer, 2012: The binning of metagenomic
257 contigs for microbial physiology of mixed cultures. *Front. Microbiol.*, **3**, 410,
258 doi:10.3389/fmicb.2012.00410.
- 259 Su, X., J. Xu, and K. Ning, 2012: Parallel-META: efficient metagenomic data analysis based on
260 high-performance computation. *BMC Syst. Biol.*, **6**, S16, doi:10.1186/1752-0509-6-S1-
261 S16.
- 262 Teeling, H., J. Waldmann, T. Lombardot, M. Bauer, and F. O. Glöckner, 2004: TETRA: a web-
263 service and a stand-alone program for the analysis and comparison of tetranucleotide
264 usage patterns in DNA sequences. *BMC Bioinformatics*, **5**, 163, doi:10.1186/1471-2105-
265 5-163.
- 266 Trindade-Silva, A. E., and Coauthors, 2012: Taxonomic and Functional Microbial Signatures of
267 the Endemic Marine Sponge Arenosclera brasiliensis. *PLoS ONE*, **7**, e39905,
268 doi:10.1371/journal.pone.0039905.
- 269 Whitman, W. B., D. C. Coleman, and W. J. Wiebe, 1998: Prokaryotes: The unseen majority. *Proc.*
270 *Natl. Acad. Sci.*, **95**, 6578–6583.
- 271 Zhang, Z., S. Schwartz, L. Wagner, and W. Miller, 2000: A greedy algorithm for aligning DNA
272 sequences. *J. Comput. Biol. J. Comput. Mol. Cell Biol.*, **7**, 203–214,
273 doi:10.1089/10665270050081478.
- 274

Figure 1

Figure 1

Workflow of the FOCUS program.

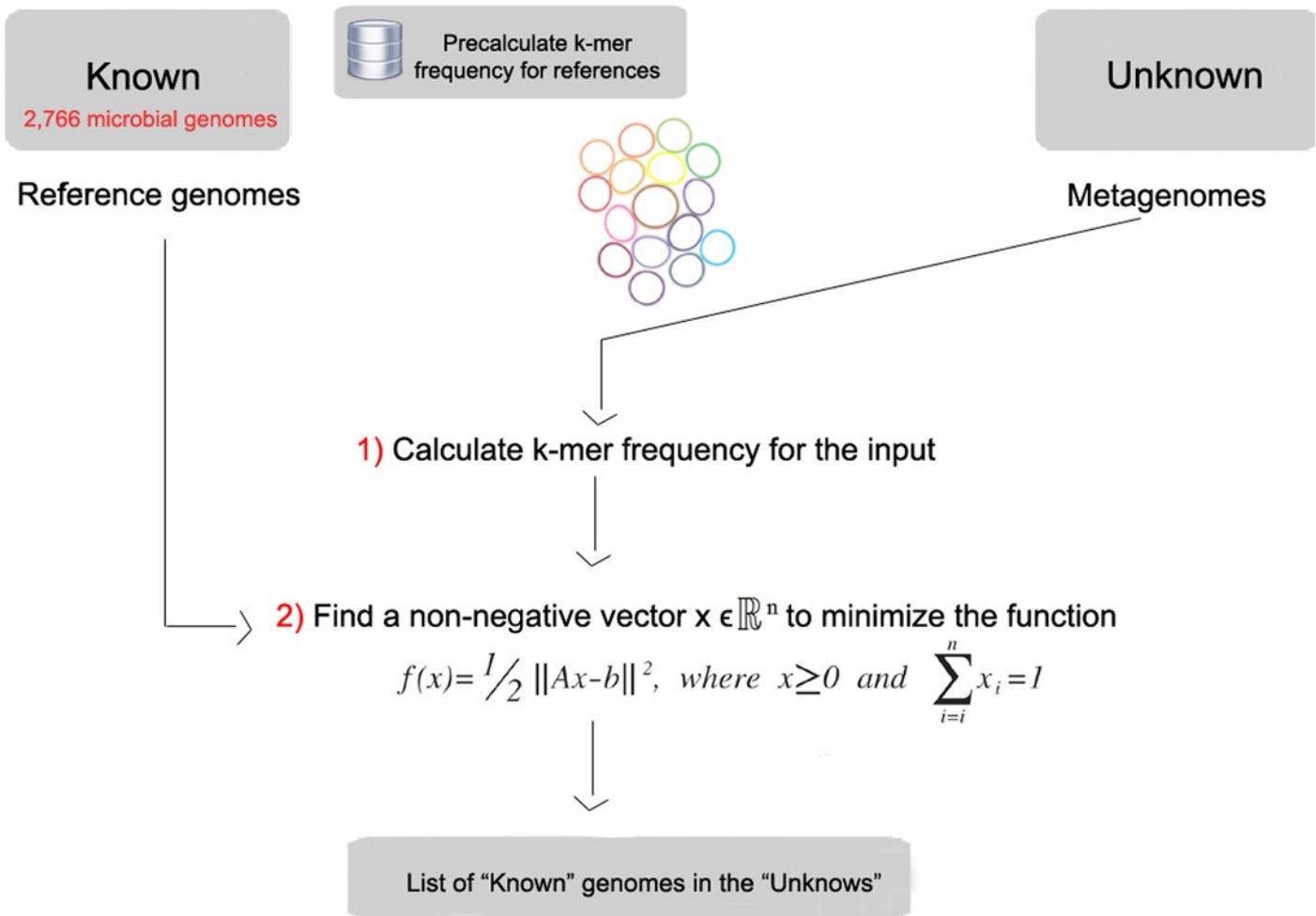


Figure 2

Figure 2

Genera-level taxonomy classification sorted by FOCUS prediction for the metagenome from a diseased human oral cavity using FOCUS, MetaPhlAn, MG-RAST, PhymBL, RAphy, Taxy, and FOCUS (mean). Error bars represent the standard deviation uncertainty in tested metagenome.

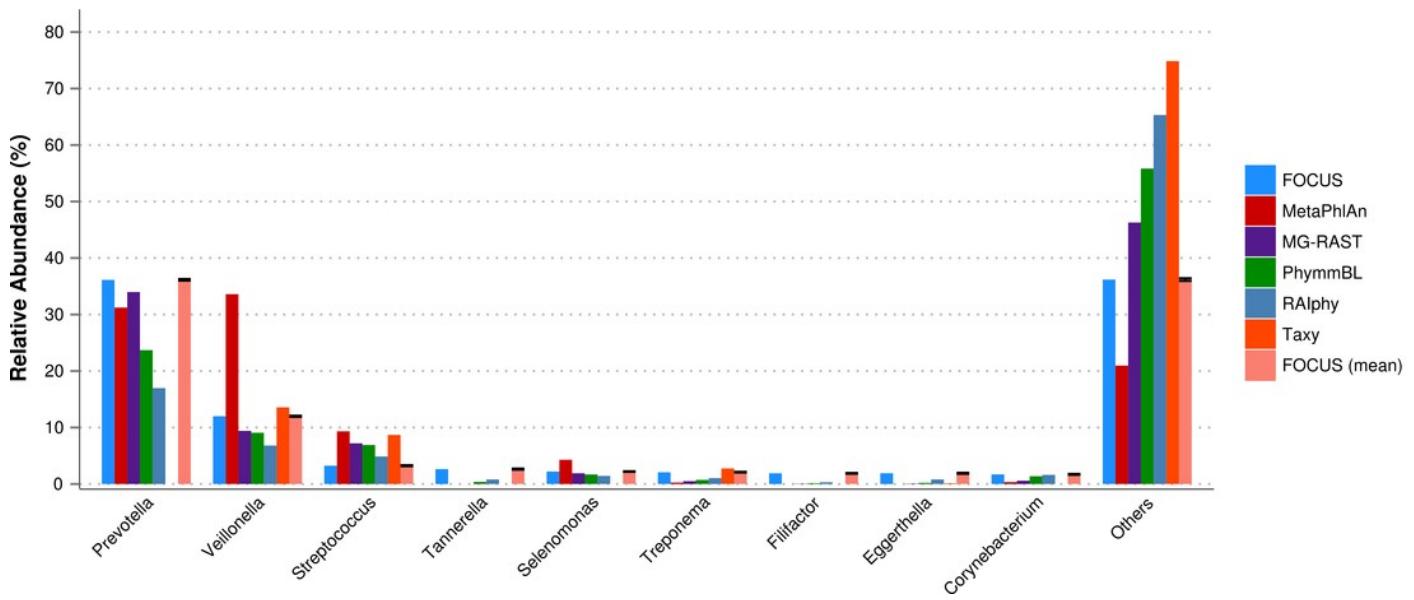


Figure 3

Figure 3

Scalability test using different sub-sets of the human oral cavity under disease metagenome using FOCUS, MetaPhlAn, MG-RAST, PhymBL, RAphy, Taxy.

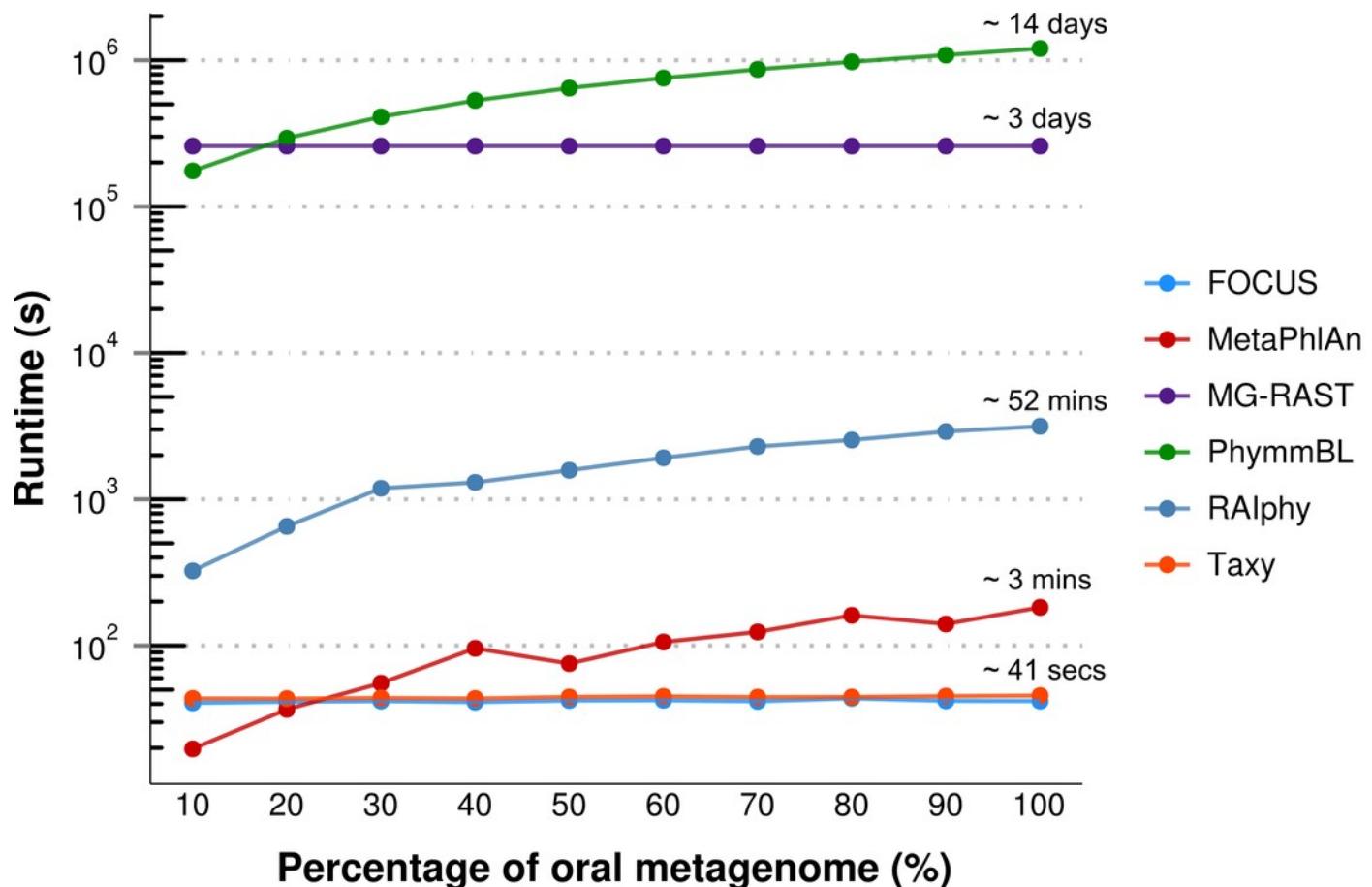


Figure 4

Figure 4

Genera-level taxonomy classification sorted by FOCUS prediction for the metagenome from a health human oral cavity using FOCUS, MetaPhlAn, MG-RAST, PhymBL, RAphy, Taxy, and FOCUS (mean). Error bars show the standard deviation uncertainty for the real metagenome.

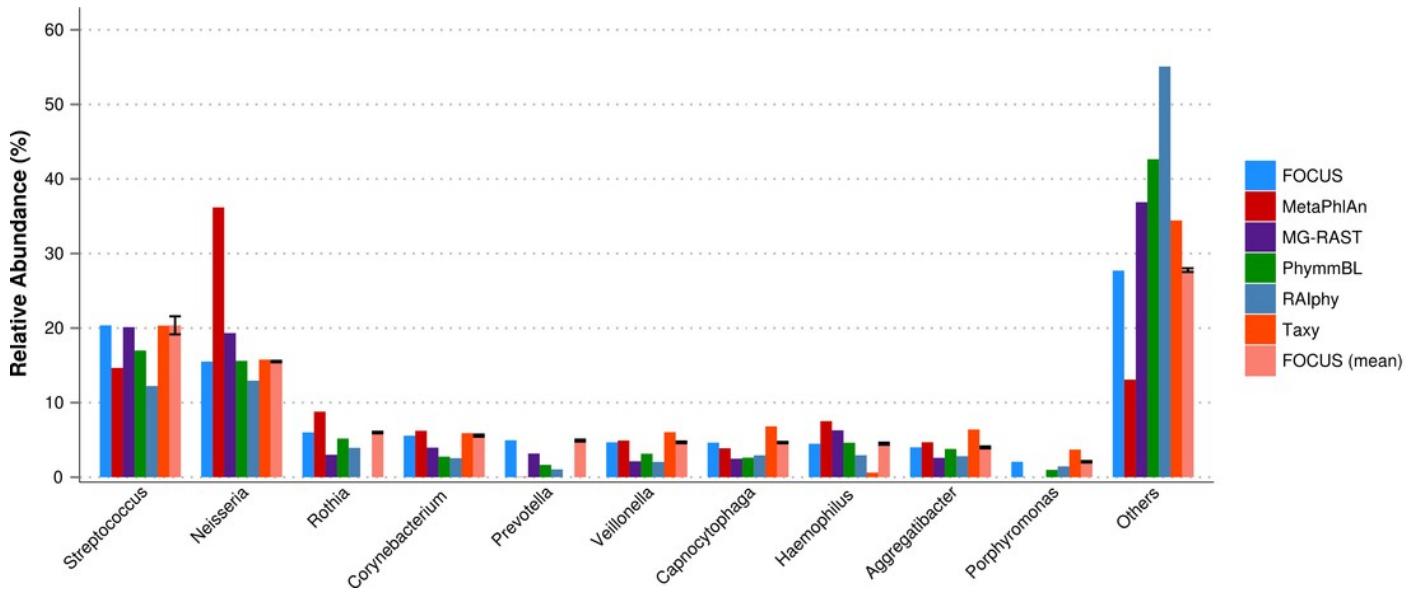


Figure 5

Figure 5

Genera-level taxonomy classification for the SimShort dataset using FOCUS, PhymBL, RAlphy, and FOCUS (mean).

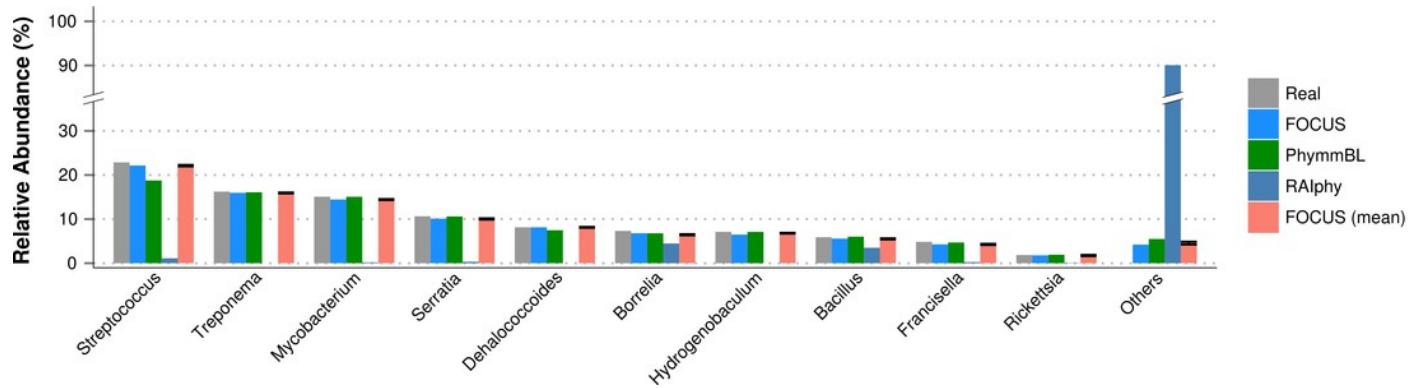


Figure 6

Figure 6

Class-level taxonomy classification for the SimHC dataset using FOCUS, PhymBL, RAlphy, and FOCUS (mean).

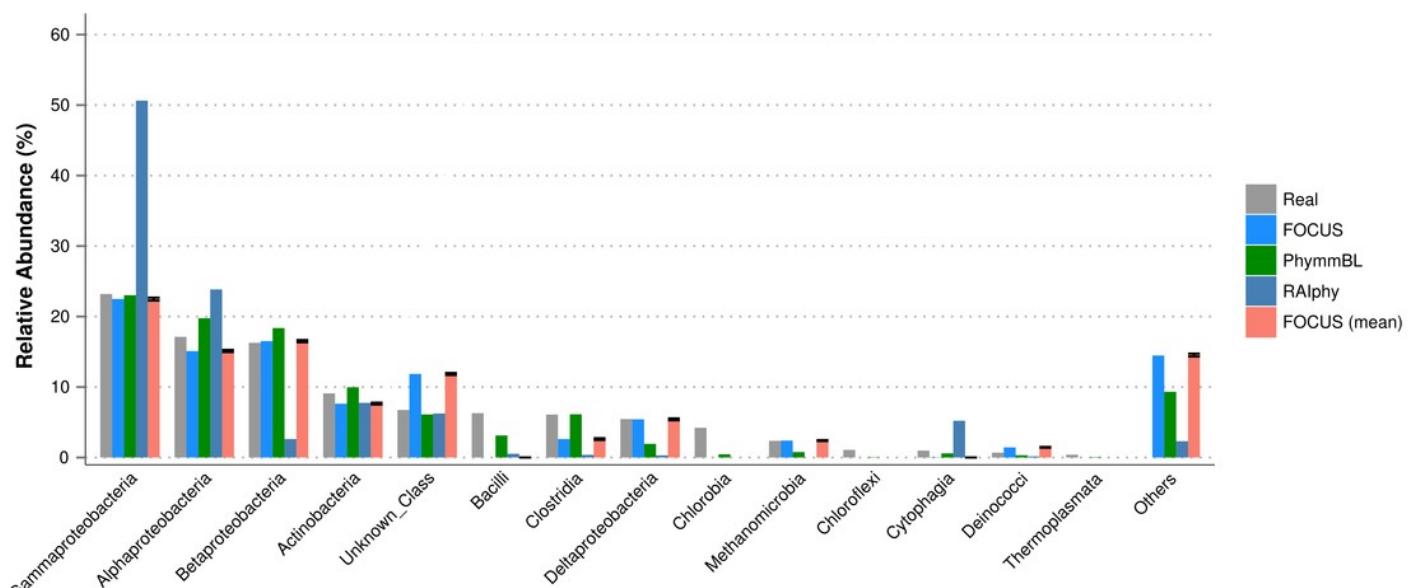


Figure 7

Figure 7

Genera-level taxonomy classification for the SimHC dataset using FOCUS, MetaPhiAn, MG-RAST, PhymnBL, RAlphy, Taxy, GASiC, and FOCUS (mean).

