# The enveomics collection: a toolbox for specialized analyses of microbial genomes and metagenomes

Luis M Rodriguez-R, Konstantinos T Konstantinidis

Genomic and metagenomic analyses are increasingly becoming commonplace in several areas of biological research, but recurrent specialized analyses are frequently reported as in-house scripts rarely available after publication. We describe the enveomics collection, a growing set of actively maintained scripts for several recurrent and specialized tasks in microbial genomics and metagenomics, and present a graphical user interface and several case studies. Our resource includes previously described as well as new algorithms such as Transformed-space Resampling In Biased Sets (TRIBS), a novel method to evaluate phylogenetic under- or over-dispersion in reference sets with strong phylogenetic bias. The enveomics collection is freely available under the terms of the Artistic License 2.0 at https://github.com/lmrodriguezr/enveomics and for online analysis at http://enve-omics.ce.gatech.edu .

1  **The enveomics collection: a toolbox for specialized analyses of**
2  **microbial genomes and metagenomes**

3  Luis M Rodriguez-R[1,2,*] and Konstantinos T. Konstantinidis[1,2,3,*]

4  [1] School of Biology, [2] Center for Bioinformatics and Computational Genomics, and [3] School of Civil and
5  Environmental Engineering, Georgia Institute of Technology, Atlanta, Georgia, 30332, USA
6  * To whom correspondence should be addressed. Tel: +1 (404) 580-0541; Fax: +1 (404) 894-8266; Email:
7  lmrodriguezr@gmail.com. Correspondence may also be addressed to K.T. Konstantinidis. Tel: +1 (404) 385-4581;
8  Fax: +1 (404) 894-8266; Email: kostas@ce.gatech.edu

9  **ABSTRACT**

10  Genomic and metagenomic analyses are increasingly becoming commonplace in several areas of
11  biological research, but recurrent specialized analyses are frequently reported as in-house scripts rarely
12  available after publication. We describe the enveomics collection, a growing set of actively maintained
13  scripts and bioinformatics algorithms for several recurrent and specialized tasks in microbial genomics
14  and metagenomics, and present a graphical user interface and several case studies. Our resource
15  includes previously described (*e.g.*, Nonpareil, MyTaxa, ANI/AAI calculator) as well as new algorithms
16  such as Transformed-space Resampling In Biased Sets (TRIBS), a novel method to evaluate
17  phylogenetic under- or over-dispersion in reference sets with strong phylogenetic bias. The enveomics
18  collection is freely available under the terms of the Artistic License 2.0 at
19  https://github.com/lmrodriguezr/enveomics and for online analysis at http://enve-omics.ce.gatech.edu.

20  **INTRODUCTION**

21  Microbial genomics and metagenomics have become key components of several areas of modern
22  research including biomedicine, epidemiology, plant and animal pathology, environmental engineering
23  and science, microbial ecology, and evolutionary biology. Specialized computational analyses in these
24  areas are hence becoming commonplace for the non-expert, often resulting in the reimplementation of
25  scripts critical for understanding and reproducing the reported results, with varying levels of quality,
26  reproducibility, and availability. While the literature analysing *ad hoc* scripts is scarce, a 2004 survey on
27  the availability of URLs reported in MEDLINE found that 19% of 1,020 analysed pages were always
28  unavailable, and only 63% were always available (Wren, 2004). Moreover, we searched the manuscripts
29  available as full-text in PubMed Central with the terms "in-house script", "in-house developed script", "in-
30  house perl script" (other languages didn't return additional results), or the same terms in plural, and found
31  1,929 matching articles (as of January 05, 2016). From these, 1,654 were related to genomics or
32  metagenomics and 449 to microbial genomics or metagenomics. We further explored the latter set of
33  manuscripts, and found that only 6% provided access to the source code (26/449), with an additional 1%
34  reporting websites no longer available (3/449) or not including the reported scripts (3/449). 3% of the
35  manuscripts explicitly indicated that the scripts were available upon request (13/449), but in only 3 cases
36  the authors provided the code within two months of the request. The large majority (90%) did not provide

37   any reference, provided references to previous publications of the same group not including the source

38   code, or referenced only the programming language in which the scripts were implemented. While this

39   survey is not a systematic analysis of availability of in-house scripts, nor does it provide quality

40   assessments, the results do underscore the prevalence of a phenomenon that undermines reproducibility

41   in studies applying microbial genomics or metagenomics. On one hand, individual tools are the basis for

42   complete and reproducible methods that are reported either in manuscripts, white papers, or standard

43   operating procedures (SOPs), but the abovementioned statistics showed that the tools infrequently

44   become available. On the other hand, a suitable alternative to providing developed tools for results

45   reproducibility is to make data on each step of the analyses publicly available, but this approach is also

46   rarely adopted, with the added issues of larger file sizes and heterogeneity, making the documentation of

47   data even more challenging than documenting code. Here, we present a growing collection of actively

48   maintained scripts for several recurrent and specialized tasks in microbial genomics and metagenomics,

49   together with comprehensive documentation, a graphical user interface, and some cases of use. Our

50   collection may also constitute a reference example for other researchers in the future, and an actively

51   maintained framework that could be collaboratively expanded.

52   **IMPLEMENTATION**

53   The enveomics collection is a multi-language set of over 70 independent scripts that accomplish

54   specialized tasks in genomics and metagenomics, including code in Ruby, Perl, AWK, Bash, and R. In

55   addition, the collection features reusable libraries that automate recurrent sub-tasks. For example, the

56   enveomics_rb Ruby library includes object-oriented representations of trees, read-placement results, sets

57   of orthologous genes, and complex (non-contiguous) sequence coordinates, together with methods for

58   accessing and downloading remote data. The R code is packaged into a single library (enveomics.R) to

59   simplify its distribution.

60   **Preferred file formats**

61   Format incompatibility between data sources and analysis tools is a common problem in bioinformatics,

62   and there are several tools and libraries dedicated to the translation between format specifications (Rice,

63   Longden & Bleasby, 2000). In order to mitigate the impact of this problem, the enveomics collection has

64   been designed to support only a reduced number of formats, with a wide range of alternative variations.

65   For example, sequence files are expected to be in FastA, but the scripts in the collection always support

66   multi-FastA and can parse variations of the definition lines and colon-lead comments (Suppl. Table S1).

67   Supported formats for other data types include tabular BLAST for similarity searches (including variations

68   with additional columns and comments), JPlace for phylogenetic read placement (Matsen et al., 2012),

69   and tables in raw text with tab-delimited columns.

70   **Access to remote servers**

71  Local data sources are often insufficient and commonly out-dated. In response, we have implemented

72  several utilities to simplify the automated access to remote databases using the Representational State

73  Transfer Application Program Interfaces (RESTful APIs) of the European Bioinformatics Institute of the

74  European Molecular Biology Laboratory (EMBL-EBI), the U.S. National Center for Biotechnology

75  Information (NCBI) E-Utilities, the Kyoto Encyclopedia of Genes and Genomes (KEGG), and the M5nr

76  (Kanehisa & Goto, 2000; Sayers et al., 2009; Wilke et al., 2012; Li et al., 2015). All the scripts using these

77  modules are categorized in Annotation/database mapping, and include additional documentation such as

78  informing the user that third-party software or database is used and thus, the latter resources should be

79  cited appropriately in any resulting publications.

80  **Enveomics-GUI**

81  The documentation and parameter descriptions for all the scripts are standardized into a set of JSON files

82  that allow the dynamic creation of Graphical User Interface (GUI) forms though the enveomics-GUI

83  package, including a set of examples and reference files (Fig. 1). The package is a collection of Ruby

84  libraries, including EnveGUI that implements graphical user interaction with Shoes 4

85  (https://github.com/shoes/shoes4). The JSON files meet the definitions of the ECMA-404 standard (Ecma

86  International, 2013), but their processing (implemented in the EnveJSON library) ignores object entries

87  with "_" key, that are utilized for comments, and implements external file inclusion using the object entries

88  with "_include" key. The package is distributed as source code (requires Shoes 4 and JRuby), as a stand-

89  alone OS-independent Java Archive (JAR), and as a bundled Mac OS X application.

90  **RESULTS**

91  **Reimplementations and novel algorithms**

92  The enveomics collection aims to simplify the use of novel and previously described algorithms for the

93  analysis of community (*e.g.*, Chao1.pl, AlphaDiversi-ty.pl, Newick.autoprune.R) and population diversity

94  (*e.g.*, BlastTab.recplot2.R, RecPlot2.find_peaks.R, CharTable.classify.rb), among other tasks in microbial

95  genomics and metagenomics. Here we describe representative modules (see also Suppl. Table S1)

96  including algorithms developed by our group.

97  *Reciprocal Best Match and Average Sequence Identity.* The detection of Reciprocal Best Matches (RBMs)

98  is a reliable method for the identification of orthology (Wolf & Koonin, 2012) that has been widely used in

99  genome-aggregate metrics of genetic relatedness (Konstantinidis & Tiedje, 2005; Goris et al., 2007).

100  Although phylogenetic reconstruction remains the gold standard for orthology detection, RBM provides a

101  fast alternative for high-throughput analyses such as genome-wide scanning. The enveomics collection

102  contains utilities for the detection of RBMs (rbm.rb) and the compilation of Orthology Groups (OGs;

103  ogs.mcl.rb), as well as the estimation of Average Nucleotide Identity (ANI; ani.rb; generally suitable for

104 comparisons of genomes assigned to the same genus) and Average Amino acid Identity (AAI; aai.rb;

105 suitable for comparisons of genomes assigned to different species).

106 *Transformed-space Resampling In Biased Sets (TRIBS).* Environmental analyses often rely on pre-

107 existing reference databases as a proxy to the presence of features in query datasets. However,

108 databases seldom represent the source of the query sets uniformly, introducing sampling biases. TRIBS

109 is a novel algorithm that reduces the impact of biased sampling by uniformly resampling reference objects

110 in a transformed space generated by Multidimensional Scaling (MDS). This enables the testing of

111 differences between a dataset and a given subset for the detection of under- or over-dispersion of traits

112 (TRIBS.test.R, TRIBS.plot-test.R). The method was originally designed for the detection of phylogenetic

113 under-dispersion of traits in groups of genomes with strong phylogenetic bias (Suppl. Fig. S1;

114 TRIBS.test.R).

115 *Automated pruning of phylogenetic trees.* The enveomics collection also features a utility to automatically

116 prune trees keeping clade representatives (Newick.autoprune.R), a useful tool for the navigation of large

117 trees such as those produced from 16S rRNA gene databases. This script iteratively extracts the

118 cophenetic matrix from a tree and removes terminal nodes with at least one other node closer than a

119 target minimum distance (by default, the first quartile of all the paired distances in the initial tree). In some

120 cases, the complete cophenetic matrix is prohibitively expensive to estimate (in the initial iterations for

121 large trees); in those cases the script takes a random sample of terminal nodes and removes sister nodes

122 (or their children) closer than the target distance. An example of a pruned tree is presented in Fig. 2B-C.

123 **Case studies using the enveomics collection**

124 *Core genome phylogenies.* Whole-genome phylogenetic reconstruction is a powerful method for the

125 resolution of evolutionary relationships. The enveomics collection includes utilities to download genomes

126 of a given species, detect RBMs between pairs of genomes, identify OGs, and identify the genes shared

127 among all the genomes in the collection –the core genome– (RefSeq.download.bash, rbm.rb, ogs.mcl.rb,

128 ogs.extract.rb). After computing independent alignments of each core OG, a concatenated alignment can

129 be generated with the options of excluding invariable sites and keeping track of coordinates (Aln.cat.rb) to

130 generate robust phylogenies with OG-specific models. In addition, the OGs can be used to estimate

131 several gene-content properties (ogs.stats.rb) and the rarefied core and pan-genomes (ogs.core-pan.rb)

132 of the species. As a less expensive alternative to the entire core genome phylogeny, one could also

133 identify and analyse only the collection of 111 single-copy genes typically present in archaeal (often ~26

134 genes) and bacterial (often ~106 genes) genomes (HMM.essential.rb). We implemented a workflow using

135 the enveomics collection, together with the multiple alignment tool Clustal Omega (Sievers et al., 2011)

136 and the phylogenetic reconstruction tool RAxML (Stamatakis, 2014) and applied it to the 17 publicly

137 available complete genomes of *Xanthomonas oryzae* (Fig. 2). The resulting phylogeny identifies known

138 pathovars and the overall structure is consistent with a previous phylogenomic reconstruction (Rodriguez-

139    R et al., 2012). The complete analysis is fully automated, and the code is deposited in the enveomics

140    collection at Examples/essential-phylogeny.bash. The execution took 31.2 minutes using two 2.9 GHz

141    processors.

142    *Gene variants in a metagenome.* Characterizing the allelic diversity of genes in metagenomes allows

143    targeted analyses of specific traits and the exploration of population discreteness and intra-population

144    variations, independent of cultivation and amplification (Caro-Quintero & Konstantinidis, 2012; Rodriguez-

145    R & Konstantinidis, 2014a). We explored the intra-population diversity of a metagenomic-recovered bin

146    (LL-70.1) using the mapping of metagenomic reads (LL_1101B; SRR948448 (Tsementzi et al., 2014))

147    from a water sample in January 2011 at Lake Lanier (GA, USA). Read mapping was performed with

148    BLAST (Altschul et al., 1990), and results were analysed and visualized using BlastTab.catsbj.pl and

149    BlastTab.recplot.R (Fig. 3), revealing small gene-content variations (panels 2 and 4), but a large allelic

150    variation and the presence of closely related organisms at about 90% ANI (panels 1 and 3). However, a

151    clear genetic discontinuity exists separating this species, as evidenced by the gap around 95% identity, a

152    phenomenon further discussed in (Caro-Quintero & Konstantinidis, 2012; Rodriguez-R & Konstantinidis,

153    2014a). The enveomics collection also includes utilities for the normalization (BlastTab.topHits_sorted.rb,

154    BlastTab.sumPerHit.pl, BlastTab.seqdepth_ZIP.pl), characterization (Chao1.pl, AlphaDiversity.pl,

155    TRIBS.test.R), and visualization (Table.barplot.R, TRIBS.plot-test.R, BlastTab.recplot2.R) of reference

156    allele distributions in a metagenome using read mapping. Additionally, the allelic diversity of a particular

157    gene of interest can be explored beyond known variants using phylogenetic read placement (Matsen,

158    Kodner & Armbrust, 2010; Berger, Krompass & Stamatakis, 2011), that can be visualized in the

159    interactive Tree of Life (iToL) (Letunic & Bork, 2007), and further explored to characterize distances to

160    known variants or ancestral nodes (JPlace.to_iToL.rb, JPlace.distances.rb) as in (Rodriguez-R et al.,

161    2015).

162    **Availability**

163    The source code for all the scripts and additional documentation are deposited and maintained at

164    https://github.com/lmrodriguezr/enveomics. The enveomics-GUI is maintained at

165    https://github.com/lmrodriguezr/enveomics-gui. In addition, we have made available a server with online

166    interfaces for select tools at http://enve-omics.ce.gatech.edu/, including the ANI and AAI calculators, and

167    previously reported tools like Nonpareil (Rodriguez-R & Konstantinidis, 2014b), a tool to estimate the level

168    of coverage in metagenomic samples, and MyTaxa (Luo, Rodriguez-R & Konstantinidis, 2014), a

169    taxonomic classification tool for sequence fragments.

170    **DISCUSSION**

171    The enveomics collection offers a wide array of tools implementing specialized recurrent tasks in

172    microbial genomics and metagenomics and is aimed for users with or without expertise in bioinformatics.

173    The collection features **(i)** a web-based interface for select tools and the complete documentation of all

174  the tools, **(ii)** a comprehensive graphical user interface (GUI), **(iii)** a command-line interface (CLI) that

175  allows integration with development platforms and automation, and **(iv)** Ruby and R application interfaces

176  (API) for developers. In addition, the collection has a language-agnostic design, allowing the

177  implementation of different tools in the most convenient language depending on available libraries or

178  other considerations. To allow this heterogeneity, all the tools are integrated using a standardized JSON-

179  based documentation scheme, allowing the incorporation of additional tools into the collection for the

180  different interfaces. Finally, examples of input data and parameters are provided to encourage the quick

181  use of the tools without dauntingly extensive user manuals.

182  A few of the scripts in our collection, in particular those implementing the most simple tasks, are

183  overlapping with those developed by others (*e.g.*, (Rice, Longden & Bleasby, 2000; Stajich et al., 2002;

184  Cock et al., 2009)). Our goal here was not to perform exhaustive comparisons to previously published

185  scripts. As explained above, these scripts were frequently not available for comparisons. Rather, the goal

186  was to put together a resource that offers easy-to-use tools for the non-bioinformatician and is

187  comprehensive with respect to recurrent tasks in microbiome research. As such, we hope that the

188  scientific community will find this resource useful, and will provide feedback on the scripts and algorithms,

189  and suggestions for further improvements.

190  **SUPPLEMENTARY DATA**

191  Supplementary data are available at NAR Online.

201  **REFERENCES**

202  Altschul SF., Gish W., Miller W., Myers EW., Lipman DJ. 1990. Basic local alignment search tool. *Journal*

203       *of molecular biology* 215:403–410. DOI: 10.1016/S0022-2836(05)80360-2.

204    Berger SA., Krompass D., Stamatakis A. 2011. Performance, Accuracy, and Web Server for Evolutionary

205          Placement of Short Sequence Reads under Maximum Likelihood. *Systematic Biology* 60:291–302.

206          DOI: 10.1093/sysbio/syr010.

207    Caro-Quintero A., Konstantinidis KT. 2012. Bacterial species may exist, metagenomics reveal.

208          *Environmental microbiology* 14:347–355. DOI: 10.1111/j.1462-2920.2011.02668.x.

209    Cock PJA., Antao T., Chang JT., Chapman BA., Cox CJ., Dalke A., Friedberg I., Hamelryck T., Kauff F.,

210          Wilczynski B., Hoon MJL de. 2009. Biopython: freely available Python tools for computational

211          molecular biology and bioinformatics. *Bioinformatics* 25:1422–1423. DOI:

212          10.1093/bioinformatics/btp163.

213    Ecma International 2013. The JSON Data Interchange Format.

214    Goris J., Konstantinidis KT., Klappenbach JA., Coenye T., Vandamme P., Tiedje JM. 2007. DNA-DNA

215          hybridization values and their relationship to whole-genome sequence similarities. *International*

216          *Journal of Systematic and Evolutionary Microbiology* 57:81–91. DOI: 10.1099/ijs.0.64483-0.

217    Kanehisa M., Goto S. 2000. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Research*

218          28:27–30. DOI: 10.1093/nar/28.1.27.

219    Konstantinidis KT., Tiedje JM. 2005. Genomic insights that advance the species definition for prokaryotes.

220          *Proceedings of the National Academy of Sciences of the United States of America* 102:2567–

221          2572. DOI: 10.1073/pnas.0409727102.

222    Letunic I., Bork P. 2007. Interactive Tree Of Life (iTOL): an online tool for phylogenetic tree display and

223          annotation. *Bioinformatics* 23:127–128. DOI: 10.1093/bioinformatics/btl529.

224    Li W., Cowley A., Uludag M., Gur T., McWilliam H., Squizzato S., Park YM., Buso N., Lopez R. 2015. The

225          EMBL-EBI bioinformatics web and programmatic tools framework. *Nucleic Acids Research*

226          43:W580–W584. DOI: 10.1093/nar/gkv279.

227   Luo C., Rodriguez-R LM., Konstantinidis KT. 2014. MyTaxa: an advanced taxonomic classifier for genomic

228        and metagenomic sequences. *Nucleic Acids Research* 42:e73–e73. DOI: 10.1093/nar/gku169.

229   Matsen FA., Hoffman NG., Gallagher A., Stamatakis A. 2012. A Format for Phylogenetic Placements. *PLoS*

230        *ONE* 7:e31009. DOI: 10.1371/journal.pone.0031009.

231   Matsen FA., Kodner RB., Armbrust EV. 2010. pplacer: linear time maximum-likelihood and Bayesian

232        phylogenetic placement of sequences onto a fixed reference tree. *BMC Bioinformatics* 11:538.

233        DOI: 10.1186/1471-2105-11-538.

234   Rice P., Longden I., Bleasby A. 2000. EMBOSS: The European Molecular Biology Open Software Suite.

235        *Trends in Genetics* 16:276–277. DOI: 10.1016/S0168-9525(00)02024-2.

236   Rodriguez-R LM., Grajales A., Arrieta-Ortiz M., Salazar C., Restrepo S., Bernal A. 2012. Genomes-based

237        phylogeny of the genus Xanthomonas. *BMC Microbiology* 12:43. DOI: 10.1186/1471-2180-12-43.

238   Rodriguez-R LM., Overholt WA., Hagan C., Huettel M., Kostka JE., Konstantinidis KT. 2015. Microbial

239        community successional patterns in beach sands impacted by the Deepwater Horizon oil spill.

240        *The ISME Journal*. DOI: 10.1038/ismej.2015.5.

241   Rodriguez-R LM., Konstantinidis KT. 2014a. Bypassing Cultivation To Identify Bacterial Species. *Microbe*

242        9:111–118.

243   Rodriguez-R LM., Konstantinidis KT. 2014b. Nonpareil: a redundancy-based approach to assess the level

244        of coverage in metagenomic datasets. *Bioinformatics* 30:629–635. DOI:

245        10.1093/bioinformatics/btt584.

246   Sayers EW., Barrett T., Benson DA., Bryant SH., Canese K., Chetvernin V., Church DM., DiCuccio M., Edgar

247        R., Federhen S., Feolo M., Geer LY., Helmberg W., Kapustin Y., Landsman D., Lipman DJ., Madden

248        TL., Maglott DR., Miller V., Mizrachi I., Ostell J., Pruitt KD., Schuler GD., Sequeira E., Sherry ST.,

249        Shumway M., Sirotkin K., Souvorov A., Starchenko G., Tatusova TA., Wagner L., Yaschenko E., Ye

250    J. 2009. Database resources of the National Center for Biotechnology Information. *Nucleic Acids*

251    *Research* 37:D5–15. DOI: 10.1093/nar/gkn741.

252 Sievers F., Wilm A., Dineen D., Gibson TJ., Karplus K., Li W., Lopez R., McWilliam H., Remmert M., Söding

253    J., Thompson JD., Higgins DG. 2011. Fast, scalable generation of high-quality protein multiple

254    sequence alignments using Clustal Omega. *Molecular Systems Biology* 7. DOI:

255    10.1038/msb.2011.75.

256 Stajich JE., Block D., Boulez K., Brenner SE., Chervitz SA., Dagdigian C., Fuellen G., Gilbert JGR., Korf I.,

257    Lapp H., Lehväslaiho H., Matsalla C., Mungall CJ., Osborne BI., Pocock MR., Schattner P., Senger

258    M., Stein LD., Stupka E., Wilkinson MD., Birney E. 2002. The Bioperl Toolkit: Perl Modules for the

259    Life Sciences. *Genome Research* 12:1611–1618. DOI: 10.1101/gr.361602.

260 Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large

261    phylogenies. *Bioinformatics* 30:1312–1313. DOI: 10.1093/bioinformatics/btu033.

262 Tsementzi D., Poretsky R., Rodriguez-R LM., Luo C., Konstantinidis KT. 2014. Evaluation of

263    metatranscriptomic protocols and application to the study of freshwater microbial communities.

264    *Environmental Microbiology Reports* 6:640–655. DOI: 10.1111/1758-2229.12180.

265 Wilke A., Harrison T., Wilkening J., Field D., Glass EM., Kyrpides N., Mavrommatis K., Meyer F. 2012. The

266    M5nr: a novel non-redundant database containing protein sequences and annotations from

267    multiple sources and associated tools. *BMC Bioinformatics* 13:141. DOI: 10.1186/1471-2105-13-

268    141.

269 Wolf YI., Koonin EV. 2012. A Tight Link between Orthologs and Bidirectional Best Hits in Bacterial and

270    Archaeal Genomes. *Genome Biology and Evolution* 4:1286–1294. DOI: 10.1093/gbe/evs100.

271 Wren JD. 2004. 404 not found: the stability and persistence of URLs published in MEDLINE.

272    *Bioinformatics* 20:668–672. DOI: 10.1093/bioinformatics/btg465.

273

274

**Peer**J Preprints

# Figure 1(on next page)

Screen captures of the enveomics GUI in Mac OS X.

(A) Initial (home) screen with search bar, listing all categories and subcategories, and highlighting a few selected and randomly picked scripts. (B) Complete list of scripts per category. (C) Task form for aai.rb pre-filled with an example. (D) Result of the aai.rb analysis. All screen captures correspond to v0.1.2. Future versions may differ.
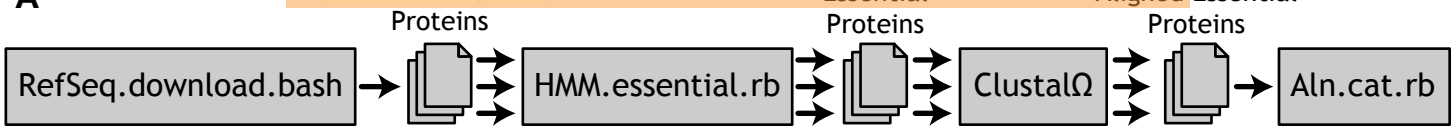
# Figure 2(on next page)

Example of a complete workflow primarily using tools from the enveomics collection applied to *Xanthomonas oryzae* genomes.

(A) The workflow uses the enveomics collection, Clustal Omega, and RAxML, to generate a phylogenetic tree based on the concatenated alignment of 105 single-copy essential genes. (B) In the resulting phylogeny, two clades form naturally corresponding to the pathovars oryzae (*Xoo*, left) and oryzicola (*Xoc*, right). Note that the tree is un-rooted, but the rooting point is suggested (vertex) based on phylogenomics of the genus (Rodriguez-R et al., 2012). The invariable sites were removed using Aln.cat.rb (35,386 sites removed) and the phylogeny was reconstructed using the remaining 552 informative sites. From the 105 detected essential genes, 22 were identical across all genomes and were excluded from the analysis. (C) A simplified version of the tree was produced by automatically pruning terminal nodes at a distance lesser than 0.01, resulting in a tree with 7 genomes (out of 17) with similar structure.
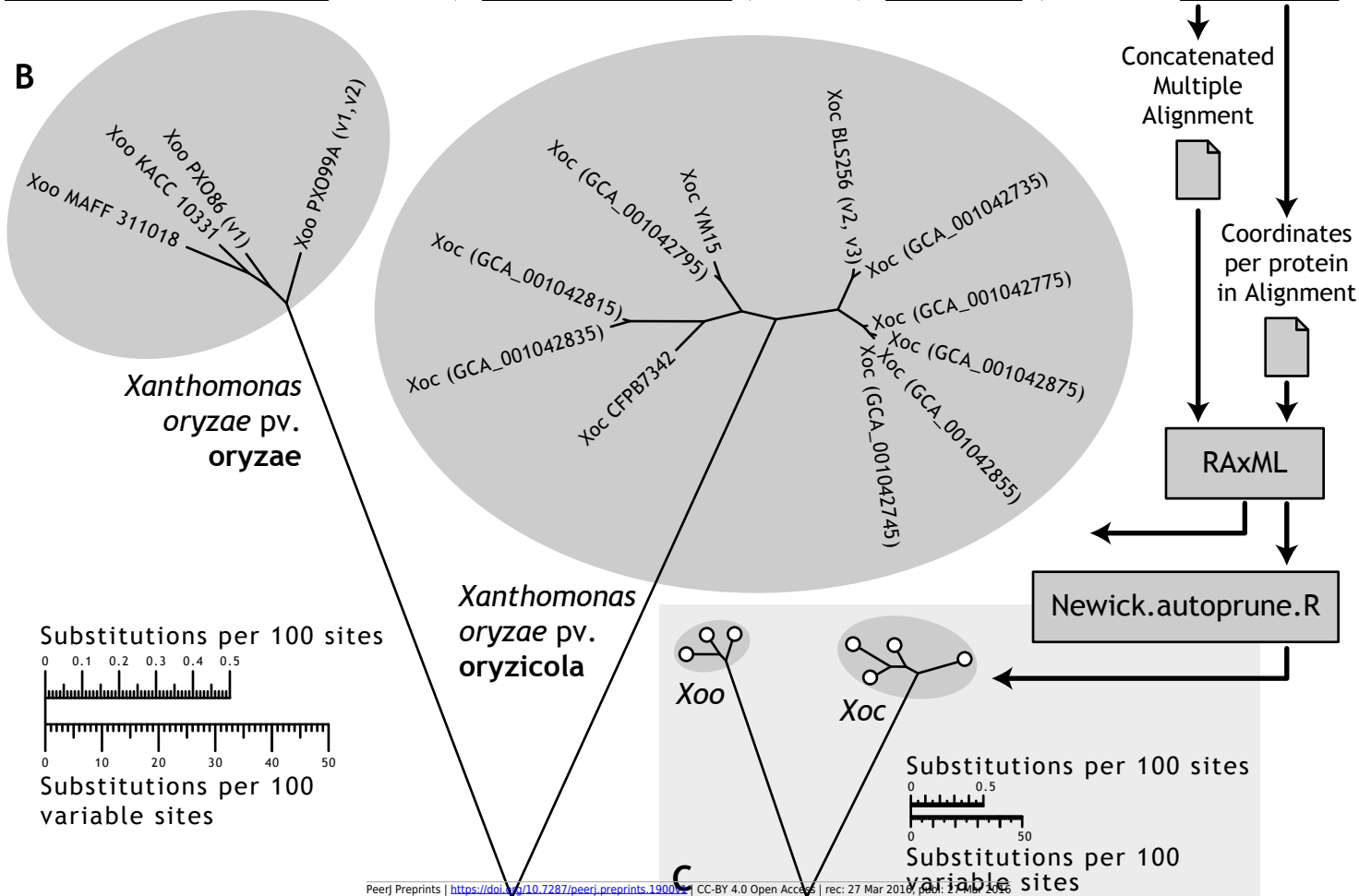
# Figure 3 (on next page)

Example of a fragment recruitment plot.

This figure showcases the result of processing a BLAST search of metagenomic short sequencing reads (150 bp long in this case; each matching read is represented by a dot in main panel 1) against a population genome sequence assembled/binned from the same metagenome (X-axis). The tabular BLAST result was parsed using BlastTab.catsbj.pl, and graphical representation was generated with the BlastTab.recplot2.R. The circled numbers 1 through 5 denote the distinct panels of the layout: **(1)** Main panel representing the reads recruited, placed by location (X-axis) and identity (Y-axis). **(2)** Sequencing depth across the reference, in logarithmic scale. Bars at the bottom represent regions without mapping reads (sequencing depth of zero). **(3)** Identity histogram of mapping reads (light gray) and smoothed spline (black), in logarithmic scale. **(4)** Sequencing depth histogram. Peaks from values above 95% identity are automatically identified as skewed normal distributions (red), with centrality measures, percentage of the reference length, and fit error (bottom-right legend) reported for each peak (marked in the right edge). **(5)** Color scale for the number of stacked reads per 2-dimensional bin in panel 1. The background of panels 1 and 3, and the line colors in panels 2 and 4, correspond to matches with identity above (dark blue) and below (light blue) a user-defined cutoff. By default, the identity cutoff is set to 95%, corresponding to the species boundary (Konstantinidis & Tiedje, 2005). See also (Rodriguez-R & Konstantinidis, 2014a) for additional discussion.