

Compact graphical representation of phylogenetic data and metadata with GraPhIAn

Francesco Asnicar, George Weingart, Timothy L Tickle, Curtis Huttenhower, Nicola Segata

The increased availability of genomic and metagenomic data poses challenges at multiple analysis levels, including visualization of very large-scale microbial and microbial community data paired with rich metadata. We developed GraPhIAn (Graphical Phylogenetic Analysis), a computational tool that produces high-quality, compact visualizations of microbial genomes and metagenomes. This includes phylogenies spanning up to thousands of taxa, annotated with metadata ranging from microbial community abundances to microbial physiology or host and environmental phenotypes. GraPhIAn has been developed as an open-source command-driven tool in order to be easily integrated into complex, publication-quality bioinformatics pipelines. It can be executed either locally or through an online Galaxy web application. We present several examples including taxonomic and phylogenetic visualization of microbial communities, metabolic functions, and biomarker discovery that illustrate GraPhIAn's potential for modern microbial and community genomics.

2 Compact graphical representation of phylogenetic data 3 and metadata with GraPhlAn

4 Francesco Asnicar ¹, George Weingart ², Timothy L Tickle ³, Curtis Huttenhower ^{2,3},
5 Nicola Segata ¹

6
7 1. Centre for Integrative Biology (CIBIO), University of Trento, Italy
8 2. Biostatistics Department, Harvard School of Public Health, USA
9 3. Broad Institute of MIT and Harvard, USA

10 Abstract

11 The increased availability of genomic and metagenomic data poses challenges at multiple
12 analysis levels, including visualization of very large-scale microbial and microbial
13 community data paired with rich metadata. We developed GraPhlAn (Graphical
14 Phylogenetic Analysis), a computational tool that produces high-quality, compact
15 visualizations of microbial genomes and metagenomes. This includes phylogenies spanning
16 up to thousands of taxa, annotated with metadata ranging from microbial community
17 abundances to microbial physiology or host and environmental phenotypes. GraPhlAn has
18 been developed as an open-source command-driven tool in order to be easily integrated
19 into complex, publication-quality bioinformatics pipelines. It can be executed either locally
20 or through an online Galaxy web application. We present several examples including
21 taxonomic and phylogenetic visualization of microbial communities, metabolic functions,
22 and biomarker discovery that illustrate GraPhlAn's potential for modern microbial and
23 community genomics.

24 Introduction

25
26 Modern high-throughput sequencing technologies provide comprehensive, large-scale
27 datasets that have enabled a variety of novel genomic and metagenomic studies. A large
28 number of statistical and computational tools have been developed specifically to tackle the
29 complexity and high-dimensionality of such datasets and to provide robust and
30 interpretable results. Visualizing data including thousands of microbial genomes or
31 metagenomes, however, remains a challenging task that is often crucial to driving
32 exploratory data mining and to compactly summarizing quantitative conclusions.

33
34 In the specific context of microbial genomics and metagenomics, next-generation
35 sequencing in particular produces datasets of unprecedented size, including thousands of
36 newly sequenced microbial genomes per month and a tremendous increase in genetic
37 diversity sampled by isolates or culture-free assays. Displaying phylogenies with thousands
38 of microbial taxa in hundreds of samples is infeasible with most available tools. This is
39 especially true when sequencing profiles need to be placed in the context of sample

40 metadata (e.g. clinical information). Among recently developed tools, iTOL (Letunic & Bork
41 2007; Letunic & Bork 2011) targets interactive analyses of large-scale phylogenies with a
42 moderate amount of overlaid metadata, whereas ETE (Huerta-Cepas et al. 2010) is a
43 Python programming toolkit focusing on tree exploration and visualization that is targeted
44 for scientific programmers, and Krona (Ondov et al. 2011) emphasizes hierarchical
45 quantitative information typically derived from metagenomic taxonomic profiles. Neither
46 of these tools provides an automatable environment for non-computationally expert users
47 in which very large phylogenies can be combined with high-dimensional metadata such as
48 microbial community abundances, host or environmental phenotypes, or microbial
49 physiological properties.

50
51 In particular, a successful high-throughput genomic visualization environment for modern
52 microbial informatics must satisfy two criteria. First, software releases must be free and
53 open-source to allow other researchers to verify and to adapt the software to their specific
54 needs and to cope with the quick evolution of data types and datasets size. Second,
55 visualization tools must be command-driven in order to be embedded in computational
56 pipelines. This allows for a higher degree of analysis reproducibility, but the software must
57 correspondingly be available for local installation and callable through a convenient
58 interface (e.g. API or general scripting language). Local installations have also the
59 advantage of avoiding the transfer of large or sensitive data to remote servers, preventing
60 potential issues with the confidentiality of unpublished biological data. Neither of these
61 criteria, of course, prevent tools from also being embeddable in web-based interfaces in
62 order to facilitate use by users with limited computational expertise (Blankenberg et al.
63 2010; Giardine et al. 2005; Goecks et al. 2010; Oinn et al. 2004), and all such tools must
64 regardless produce informative, clear, detailed, and publication-ready visualizations.

65 **Materials & Methods**

66 GraPhlAn is a new tool for compact and publication-quality representation of circular
67 taxonomic and phylogenetic trees with potentially rich sets of associated metadata. It was
68 developed primarily for microbial genomic and microbiome-related studies in which the
69 complex phylogenetic/taxonomic structure of microbial communities needs to be
70 complemented with quantitative and qualitative sample-associated metadata. GraPhlAn is
71 available at <http://cibiocm.bitbucket.org/tools/graphlan.html>.

72 **Implementation strategy**

73 GraPhlAn is composed by two Python modules: one for drawing the image and one for
74 adding annotations to the tree. GraPhlAn exploits the annotation file to highlight and
75 personalize the appearance of the tree and of the associated information. The annotation
76 file does not perform any modifications to the structure of the tree, but it just changes the
77 way in which nodes and branches are displayed. Internally, GraPhlAn uses the matplotlib
78 library (Hunter 2007) to perform the drawing functions.

79 **The export2graphlan module**

80 Export2graphlan is a framework to easily integrate GraPhlAn into already existing
81 bioinformatics pipelines. Export2graphlan makes use of two external libraries: the pandas

82 python library (McKinney 2012) and the BIOM library, only when BIOM files are given as
83 input.
84 Export2graphlan can take as input two files: the result of the analysis of MetaPhlAn (either
85 version 1 or 2) or HUMAnN, and the result of the analysis of LEfSe. At least one of these two
86 input files is mandatory. Export2graphlan will then produce a tree file and an annotation
87 file that can be used with GraPhlAn. In addition, export2graphlan can take as input a BIOM
88 file (either version 1 or 2).
89 Export2graphlan performs an analysis on the abundance values and, if present, on the LDA
90 score assigned by LEfSe, to annotate and highlight the most abundant clades and the ones
91 found to be biomarkers. Through a number of parameters the user can control the
92 annotations produced by export2graphlan.

93 Results and Discussion

94 Plotting taxonomic trees with clade annotations

95 The simplest structures visualizable by GraPhlAn include taxonomic trees (i.e. those
96 without variable branch lengths) with simple clade or taxon nomenclature labels. These
97 can be combined with quantitative information such as taxon abundances, phenotypes, or
98 genomic properties. GraPhlAn provides separate visualization options for trees (thus
99 potentially unannotated) and their annotations, the latter of which (the annotation
100 module) attaches metadata properties using the PhyloXML format (Han & Zmasek 2009).
101 This annotation and subsequent metadata visualization process (**Fig. 1**) can be repeatedly
102 applied to the same tree.

103
104 The GraPhlAn tree visualization (plotting module) takes as input a tree represented in any
105 one of the most common data formats: Newick, Nexus (Maddison et al. 1997), PhyloXML
106 (Han & Zmasek 2009), or plain text. Without annotations, the plotting module generates a
107 simple version of the tree (**Fig. 1A**), but the process can then continue by adding a diverse
108 set of visualization annotations. Annotations can affect the appearance of the tree at
109 different levels, including its global appearance (“global options” e.g. the size of the image,
110 **Fig. 1B**), the properties of subsets of nodes and branches (“node options” e.g. the color of a
111 taxon, **Fig. 1C**), and the background features used to highlight sub-trees (“label options” e.g.
112 the name of a species containing multiple taxa, **Fig. 1D**). A subset of the available
113 configurable options includes the thickness of tree branches, their colors, highlighting
114 background colors and labels of specific sub-trees, and the sizes and shapes of individual
115 nodes. Wild cards are supported to share graphical and annotation details among sub-trees
116 by affecting all the descendants of a clade or its terminal nodes only. These features in
117 combination aim to conveniently highlight specific sub-trees and metadata patterns of
118 interest.

119
120 Additional taxon-specific features can be plotted as so-called external rings when not
121 directly embedded into the tree. External rings are drawn just outside the area of the tree
122 and can be used to display specific information about leaf taxa, such as abundances of each
123 species in different conditions/environments or their genome sizes. The shapes and forms
124 of these rings are also configurable; for example, in **Fig. 1E** (“set external ring options”), the

125 elements of the innermost external ring are triangular, indicating the directional sign of a
126 genomic property. The second, third, and fourth external rings show leaf-specific features,
127 using a heatmap gradient from blank to full color. Finally, the last external ring is a bar-plot
128 representing a continuous property of leaf nodes of the tree.

129 **Compact representations of phylogenetic trees with associated metadata**

130 Visualizing phylogenetic structures and their relation to external metadata is particularly
131 challenging when the dimension of the internal structure is large. Mainly as a consequence
132 of the low cost of sequencing, current research in microbial genomics and metagenomics
133 needs indeed to visualize a considerable amount of phylogenetic data. GraPhlAn can easily
134 handle such cases, as illustrated here in an example of a large phylogenetic tree (3,737 taxa,
135 provided as a PhyloXML file in the software repository, see Availability section) with
136 multiple types of associated metadata (**Fig. 2**).

137

138 Specifically, we used GraPhlAn to display the microbial tree of life as inferred by
139 PhyloPhlAn (Segata et al. 2013), annotating this evolutionary information with genome-
140 specific metadata (**Fig. 2**). In particular, we annotated the genome contents related to
141 seven functional modules from the KEGG database (Kanehisa et al. 2011), specifically two
142 different ATP synthesis machineries (M00157: F-type ATPase and M00159: V/A-type
143 ATPase) and five modules for bacterial fatty acid metabolism (M00082: Fatty acid
144 biosynthesis, initiation, M00083: Fatty acid biosynthesis elongation, M00086: acyl-CoA
145 synthesis, M00087: beta-Oxidation, and M00088: Ketone body biosynthesis). We then also
146 annotated genome size as an external circular bar plot.

147

148 As expected, it is immediately visually apparent that the two types of ATPase are almost
149 mutually exclusive within available genome annotations, with the V/A-type ATPase
150 (module M00159) present mainly in *Archaea* and the F-type ATPase (module M000157)
151 mostly characterizing *Bacteria*. Some exceptions are easily identifiable: *Thermi* and
152 *Clamydophilia*, for instance, completely lack the F-type ATPase, presenting only the
153 typically archaea-specific V/A-type ATPase. As previously discussed in the literature (Cross
154 & Müller 2004; Mulkidjanian et al. 2007), this may due to the acquisition of V/A-type
155 ATPase by horizontal gene transfer and the subsequent loss of the F-type ATPase
156 capability. Interestingly, some species such as those in the *Streptococcus* genus and some
157 *Clostridia* still show both ATPase systems in their genomes.

158

159 With respect to fatty acid metabolism, some clades - including organisms such as
160 *Mycoplasmas* - completely lack any of the targeted pathways. Indeed, *Mycoplasmas* are the
161 smallest living cells yet discovered, lacking a cell wall (Razin 1992) and demonstrating an
162 obligate parasitic lifestyle. Since they primarily exploit host molecular capabilities,
163 *Mycoplasmas* do not need to be able to fulfill all typical cell functions, and this is also
164 indicated by the plotted very short genome sizes. *Escherichia*, on the other hand, has a
165 much longer genome, and all the considered fatty acid metabolism capabilities are present.
166 These evolutionary aspects are well known in the literature, GraPhlAn permits them and
167 other phylogeny-wide genomic patterns to be easily visualized for further hypothesis
168 generation.

169 Visualizing microbiome biomarkers

170 GraPhlAn provides a means for displaying either phylogenetic (trees with branch lengths)
171 or taxonomic (trees without branch length) data generated by other metagenomic analysis
172 tools. For instance, we show here examples of GraPhlAn plots for taxonomic profiles (**Fig.**
173 **3**), functional profiles (**Fig. 4**), and specific features identified as biomarkers (**Fig. 3** and **4**).
174 In these plots, GraPhlAn highlights microbial sub-trees that are found to be significantly
175 differentially abundant by LEfSe (Segata et al. 2011), along with their effect sizes as
176 estimated by linear discriminant analysis (LDA). To enhance biomarker visualization, we
177 annotated them in the tree with a shaded background color and with clade names as labels,
178 with decreasing font sizes for internal levels. To represent the effect size, we scaled the
179 node color from black (low LDA score) to full color (high LDA score).

180

181 **Fig. 3** shows the taxonomic tree of biomarkers (significantly differential clades) resulting
182 from a contrast gut metagenome profiles from the Human Microbiome Project (HMP)
183 (Huttenhower et al. 2012) and MetaHIT samples (Qin et al. 2010). Only samples from
184 healthy individuals in the latter cohort were included. The filtered dataset was analyzed
185 using LEfSe (Segata et al. 2011) and the cladogram obtained using the *export2graphlan*
186 script provided with GraPhlAn and discussed in the following section. As expected, the
187 image highlights that *Firmicutes* and *Bacteroides* are the two most abundant taxa in the
188 healthy gut microbiome (David et al. 2014; Wu et al. 2011). The *Bacteroidetes* phylum
189 contains many clades enriched in the HMP dataset, while *Firmicutes* show higher
190 abundances for MetaHIT samples. GraPhlAn can thus serve as a visual tool for inspecting
191 specific significant differences between conditions or cohorts.

192

193 Functional ontologies can be represented by GraPhlAn in a similar way and provide
194 complementary features to the types of taxonomic analyses shown above. Metabolic
195 profiles quantified by HUMAnN (Abubucker et al. 2012) using KEGG (Kanehisa et al. 2014)
196 from the same set of HMP and MetaHIT samples are again contrasted on multiple
197 functional levels in **Fig. 4**. The tree highlights three different broad sets of metabolic
198 pathways: Environmental Information Processing, Genetic Information Processing, and
199 Metabolism, with the last being the largest subtree. More specific metabolic functions are
200 specifically enriched in the HMP cohort, such as Glycolysis and the Citrate cycle, or in the
201 MetaHIT cohort, such as Sulfur Metabolism and Vitamin B6 Metabolism. This illustrates
202 GraPhlAn's use with different types of data, such as functional trees in addition to
203 taxonomies or phylogenies. By properly configuring input parameters of *export2graphlan*,
204 we automatically obtained both **Fig. 3** and **Fig. 4** (bash scripts used for these operations
205 are available in the GraPhlAn software repository).

206 Reproducible integration with existing analysis tools and pipelines

207 Graphical representations are usually a near- final step in the complex computational and
208 metagenomic pipelines, and automating their production is crucial for convenient but
209 reproducible analyses. To this end, GraPhlAn has been developed with command-driven
210 automation in mind, as well as flexibility in the input "annotation file" so as to be easily
211 generated by automated scripts. Depending on the specific analysis, these scripts can focus
212 on a diverse set of commands to highlight the features of interest. Despite this flexibility,
213 we further tried to ease the integration of GraPhlAn by providing automatic offline

214 conversions for some of the available metagenomic pipelines and by embedding it into the
215 well-established Galaxy web framework (Blankenberg et al. 2010; Giardine et al. 2005;
216 Goecks et al. 2010).

217

218 In order to automatically generate GraPhlAn plots from a subset of available shotgun
219 metagenomic tools comprising MetaPhlAn (for taxonomic profiling), HUMAnN (for
220 metabolic profiling), and LEfSe (for biomarker discovery), we developed a script named
221 “export2graphlan” able to convert the outputs of these tools into GraPhlAn input files as
222 schematized in **Fig. 5**. This conversion software is also meant to help biologists by
223 providing initial, automated input files for GraPhlAn that can then be manually tweaked for
224 specific needs such as highlighting clades of particular interest. The export2graphlan
225 framework can further accept the widely adopted BIOM format, both versions 1 and 2
226 (McDonald et al. 2012). This makes it possible to readily produce GraPhlAn outputs from
227 other frameworks such as QIIME (Caporaso et al. 2010) and mothur (Schloss et al. 2009)
228 for 16S rRNA sequencing studies.

229

230 A web-based deployment of the GraPhlAn application is available to the public via Galaxy at
231 <http://huttenhower.sph.harvard.edu/galaxy/>. The Galaxy interface of GraPhlAn consists of
232 four processing modules: (1) *Upload file*, that manages the upload of the input data into
233 Galaxy; (2) *GraPhlAn Annotate Tree*, which allows the user to specify the annotations that
234 will be applied to the final image; (3) *Add Rings to tree*, an optional step to select an already
235 uploaded file in Galaxy that will be used as an annotation file for the external rings; and (4)
236 *Plot tree*, that sets some image parameters such as the size, the resolution, and the output
237 format.

238 Conclusions

239 We present GraPhlAn, a new method for generating high-quality circular phylogenies
240 potentially integrated with diverse, high-dimensional metadata. We provided several
241 examples showing the application of GraPhlAn to phylogenetic, functional, and taxonomic
242 summaries. The system has already been used for a variety of additional visualization
243 tasks, including highlighting the taxonomic origins of metagenomic biomarkers (Segata et
244 al. 2012; Segata et al. 2011; Shogan et al. 2014; Xu et al. 2014), exposing specific
245 microbiome metabolic enrichments within a functional ontology (Abubucker et al. 2012;
246 Sczesnak et al. 2011), and representing 16S rRNA sequencing results (Ramirez et al. 2014).
247 GraPhlAn is, however, not limited to microbiome data and has additionally been applied to
248 animal and plant taxonomies (Tree of Sex Consortium 2014) and to large prokaryotic
249 phylogenies built using reference genomes (Baldini et al. 2014; Chai et al. 2014; Langille et
250 al. 2013; Segata et al. 2013).

251

252 Compared to the other existing state-of-the-art approaches such as Krona (Ondov et al.
253 2011) and iTOL (Letunic & Bork 2007; Letunic & Bork 2011), GraPhlAn provides greater
254 flexibility, configuration, customization, and automation for publication reproducibility. It
255 is both easily integrable into automated computational pipelines and can be used
256 conveniently online through the Galaxy-based web interface. The software is available

257 open-source, and the features highlighted here illustrate a number of ways in which its
258 visualization capabilities can be integrated into microbial and community genomics to
259 display large tree structures and corresponding metadata.

260 Data and software availability

261 Description of the datasets and figure generation

262 The data of the taxonomic trees presented in **Fig. 1** is available in the *guide* folder, inside
263 the *examples* directory of the GraPhlAn repository
264 (<https://bitbucket.org/nsegata/graphlan>). This same image is thoroughly described under
265 the “A step-by-step example” section, in the GraPhlAn wiki included in the repository.

266
267 The genomic data used for the Tree of Life in **Fig. 2** was obtained from the Integrated
268 Microbial Genomes (IMG) data management system of the U.S. Department of Energy Joint
269 Genome Institute (DOE JGI) 2.0 dataset (http://jgi.doe.gov/news_12_1_06/). From the
270 KEGG database (Kanehisa & Goto 2000; Kanehisa et al. 2014) we focused on the following
271 modules: M00082, M00083, M00086, M00087, M00088, M00157, and M00159. The input
272 data for drawing **Fig. 2** is available in the *PhyloPhlAn* folder under the *examples* directory of
273 the GraPhlAn repository.

274
275 In **Fig. 3**, to comprehensively characterize the asymptomatic human gut microbiota, we
276 combined 224 fecal samples (>17 million reads) from the Human Microbiome Project
277 (HMP) (Human Microbiome Project 2012a; Human Microbiome Project 2012b) and the
278 MetaHIT (Qin et al. 2010) projects, two of the largest gut metagenomic collections
279 available. The taxonomic profiles were obtained by applying MetaPhlAn2. The 139 fecal
280 samples from the HMP can be accessed at <http://hmpdacc.org/HMASM/>, whereas the 85
281 fecal samples from MetaHIT were downloaded from the European Nucleotide Archive
282 (<http://www.ebi.ac.uk/ena/>, study accession number ERP000108). The input files for
283 obtaining this image with GraPhlAn are present into the *examples* folder of the repository,
284 inside the *hmp_metahit* directory. The two input files represent the merge result of the
285 MetaPhlAn analysis (*hmp_metahit.txt*) and the LEfSe result on the first file
286 (*hmp_metahit.lefse.txt*). The bash script provided exploits the export2graphlan capabilities
287 to generate the annotation file.

288
289 The functional profiles used in **Fig. 4** are the reconstruction of the metabolic activities of
290 microbiome communities. The HUMAnN pipeline (Abubucker et al. 2012) infers
291 community function directly from short metagenomic reads, using the KEGG ortholog (KO)
292 groups. HUMAnN was run on the same samples of **Fig. 3**. The dataset is available on-line at
293 <http://www.hmpdacc.org/HMMRC/>. As for the previous figure, the input files for obtaining
294 **Fig. 4** are uploaded in the *hmp_metahit_functional* folder, inside the *examples* directory of
295 the repository. The two files (*hmp_metahit_functional.txt* and
296 *hmp_metahit_functional.lefse.txt*) represent the result of HUMAnN on the HMP and MetaHIT
297 datasets and the result of LEfSe executed on the former file. The bash script provided
298 executes export2graphlan for generating the annotation file and then invoking GraPhlAn
299 for plotting the functional tree.

300

301 The dataset of supplementary **Fig. S1** refer to a 16S rRNA amplicon experiment.
302 Specifically, it consists of 454 FLX Titanium sequences spanning the V3 to V5 variable
303 regions, obtained from 24 healthy samples (12 male and 12 female) for a total of 301
304 samples. Detailed protocols used for enrollment, sampling, DNA extraction, 16S
305 amplification and sequencing are available on the Human Microbiome Project Data
306 Analysis and Coordination Center website HMP Data Analysis and Coordination Center
307 (http://www.hmpdacc.org/tools_protocols/tools_protocols.php). This data are pilot
308 samples from the HMP project (Segata et al. 2011). The input files for obtaining this image
309 is available in the *examples* folder of the export2graphlan repository
310 (<https://bitbucket.org/CibioCM/export2graphlan>), inside the *hmp_aerobiosis* directory.
311 The two files represent the taxonomic tree of the HMP project and the results of LEfSe
312 executed on the same data.

313

314 In the supplementary **Fig. S2** we used the saliva microbiome profiles obtained by 16S rRNA
315 sequencing on the IonTorrent platform (amplifying the hypervariable region V3). The
316 dataset comprises a total of 13 saliva samples from healthy subjects as described in (Dassi
317 et al. 2014) and it is available in the NCBI Short Read Archive
318 (<http://www.ncbi.nlm.nih.gov/sra>). The input BIOM file for drawing this image is available
319 in the *saliva_microbiome* directory inside the *examples* folder of the GraPhlAn repository.

320

321 For the supplementary **Fig. S3** data represent the temporal dynamics of the human vaginal
322 microbiota, and were taken from the study of (Gajer et al. 2012). Data were obtained by
323 16S rRNA using the 454 pyrosequencing technology (sequencing the V1 and V2
324 hypervariable regions). The dataset is composed of samples from 32 women that self-
325 collected samples twice a week for 16 weeks. The input file, provided in BIOM format, is
326 present in the *vaginal_microbiota* folder inside the *examples* directory of the GraPhlAn
327 repository.

328 **Software repository, dependences, and user support**

329 GraPhlAn is freely available (<http://cibiocm.bitbucket.org/tools/graphlan.html>) and
330 released open-source in Bitbucket (<https://bitbucket.org/nsegata/graphlan>) with a set of
331 working examples and a complete tutorial that guides users throughout its functionality.
332 GraPhlAn uses the matplotlib library (Hunter 2007). GraPhlAn is also available via a public
333 Galaxy instance at <http://huttenhower.sph.harvard.edu/galaxy/>

334

335 Export2graphlan is freely available and released open-source in Bitbucket
336 (<https://bitbucket.org/CibioCM/export2graphlan>) along with a number of examples
337 helpful for testing if everything is correctly configured and installed. The export2graphlan
338 repository is also present as a sub-repository inside the GraPhlAn repository. The
339 export2graphlan module exploits the pandas library (McKinney 2012) and the BIOM
340 library (McDonald et al. 2012).

341

342 Both GraPhlAn and export2graphlan are supported through the Google group “GraPhlAn-
343 users” (<https://groups.google.com/forum/#!forum/graphlan-users>), available also as a
344 mailing list at: graphlan-users@googlegroups.com.

345 Acknowledgements

346 We would like to thank the members of the Segata and Huttenhower labs for helpful
347 suggestions, the WebValley team and participants for inspiring comments and tests, and
348 the users that tried the alpha version of GraPhlAn providing invaluable feedback to
349 improve the software.

350 References

- 351 Abubucker S, Segata N, Goll J, Schubert AM, Izard J, Cantarel BL, Rodriguez-Mueller B,
352 Zucker J, Thiagarajan M, Henrissat B, White O, Kelley ST, Methe B, Schloss PD,
353 Gevers D, Mitreva M, and Huttenhower C. 2012. Metabolic reconstruction for
354 metagenomic data and its application to the human microbiome. *PLoS Comput Biol*
355 8:e1002358.
- 356 Baldini F, Segata N, Pompon J, Marcenac P, Robert Shaw W, Dabire RK, Diabate A, Levashina
357 EA, and Catteruccia F. 2014. Evidence of natural Wolbachia infections in field
358 populations of *Anopheles gambiae*. *Nat Commun* 5:3985.
- 359 Blankenberg D, Von Kuster G, Coraor N, Ananda G, Lazarus R, Mangan M, Nekrutenko A,
360 and Taylor J. 2010. Galaxy: a web-based genome analysis tool for experimentalists.
361 *Curr Protoc Mol Biol* Chapter 19:Unit 19 10 11-21.
- 362 Caporaso J, Kuczynski J, Stombaugh J, Bittinger K, Bushman F, Costello E, Fierer N, Pena A,
363 Goodrich J, Gordon J, Huttley G, Kelley S, Knights D, Koenig J, Ley R, Lozupone C,
364 McDonald D, Muegge B, Pirrung M, Reeder J, Sevinsky J, Turnbaugh P, Walters W,
365 Widmann J, Yatsunencko T, Zaneveld J, and Knight R. 2010. QIIME allows analysis of
366 high-throughput community sequencing data. *Nat Methods* 7:335 - 336.
- 367 Chai J, Kora G, Ahn T-H, Hyatt D, and Pan C. 2014. Functional phylogenomics analysis of
368 bacteria and archaea using consistent genome annotation with UniFam. *BMC*
369 *evolutionary biology* 14:207.
- 370 Cross RL, and Müller V. 2004. The evolution of A-, F-, and V-type ATP synthases and
371 ATPases: reversals in function and changes in the H⁺/ATP coupling ratio. *FEBS*
372 *Letters* 576:1-4.
- 373 Dassi E, Ballarini A, Covello G, Quattrone A, Jousson O, De Sanctis V, Bertorelli R, Denti MA,
374 and Segata N. 2014. Enhanced microbial diversity in the saliva microbiome induced
375 by short-term probiotic intake revealed by 16S rRNA sequencing on the IonTorrent
376 PGM platform. *Journal of Biotechnology*.
- 377 David LA, Maurice CF, Carmody RN, Gootenberg DB, Button JE, Wolfe BE, Ling AV, Devlin
378 AS, Varma Y, Fischbach MA, Biddinger SB, Dutton RJ, and Turnbaugh PJ. 2014. Diet
379 rapidly and reproducibly alters the human gut microbiome. *Nature* 505:559-563.
- 380 Gajer P, Brotman RM, Bai G, Sakamoto J, Schutte UM, Zhong X, Koenig SS, Fu L, Ma ZS, Zhou
381 X, Abdo Z, Forney LJ, and Ravel J. 2012. Temporal dynamics of the human vaginal
382 microbiota. *Sci Transl Med* 4:132ra152.
- 383 Giardine B, Riemer C, Hardison RC, Burhans R, Elnitski L, Shah P, Zhang Y, Blankenberg D,
384 Albert I, Taylor J, Miller W, Kent WJ, and Nekrutenko A. 2005. Galaxy: a platform for
385 interactive large-scale genome analysis. *Genome Res* 15:1451-1455.

- 386 Goecks J, Nekrutenko A, Taylor J, and Galaxy T. 2010. Galaxy: a comprehensive approach for
387 supporting accessible, reproducible, and transparent computational research in the
388 life sciences. *Genome Biol* 11:R86.
- 389 Han MV, and Zmasek CM. 2009. phyloXML: XML for evolutionary biology and comparative
390 genomics. *BMC bioinformatics* 10:356.
- 391 Huerta-Cepas J, Dopazo J, and Gabaldon T. 2010. ETE: a python Environment for Tree
392 Exploration. *BMC bioinformatics* 11:24.
- 393 Human Microbiome Project C. 2012a. A framework for human microbiome research.
394 *Nature* 486:215-221.
- 395 Human Microbiome Project C. 2012b. Structure, function and diversity of the healthy
396 human microbiome. *Nature* 486:207-214.
- 397 Hunter JD. 2007. Matplotlib: A 2D graphics environment. *Computing in Science &*
398 *Engineering* 9:90-95.
- 399 Huttenhower C, Gevers D, Knight R, Abubucker S, Badger JH, Chinwalla AT, Creasy HH, Earl
400 AM, Fitzgerald MG, and Fulton RS. 2012. Structure, function and diversity of the
401 healthy human microbiome. *Nature* 486:207-214.
- 402 Kanehisa M, and Goto S. 2000. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic*
403 *Acids Res* 28:27-30.
- 404 Kanehisa M, Goto S, Sato Y, Furumichi M, and Tanabe M. 2011. KEGG for integration and
405 interpretation of large-scale molecular data sets. *Nucleic acids research:gkr988*.
- 406 Kanehisa M, Goto S, Sato Y, Kawashima M, Furumichi M, and Tanabe M. 2014. Data,
407 information, knowledge and principle: back to metabolism in KEGG. *Nucleic Acids*
408 *Res* 42:D199-205.
- 409 Langille MG, Zaneveld J, Caporaso JG, McDonald D, Knights D, Reyes JA, Clemente JC,
410 Burkepile DE, Vega Thurber RL, Knight R, Beiko RG, and Huttenhower C. 2013.
411 Predictive functional profiling of microbial communities using 16S rRNA marker
412 gene sequences. *Nat Biotechnol* 31:814-821.
- 413 Letunic I, and Bork P. 2007. Interactive Tree Of Life (iTOL): an online tool for phylogenetic
414 tree display and annotation. *Bioinformatics* 23:127-128.
- 415 Letunic I, and Bork P. 2011. Interactive Tree Of Life v2: online annotation and display of
416 phylogenetic trees made easy. *Nucleic Acids Res* 39:W475-478.
- 417 Maddison DR, Swofford DL, and Maddison WP. 1997. NEXUS: an extensible file format for
418 systematic information. *Systematic Biology* 46:590-621.
- 419 McDonald D, Clemente J, Kuczynski J, Rideout J, Stombaugh J, Wendel D, Wilke A, Huse S,
420 Hufnagle J, Meyer F, Knight R, and Caporaso J. 2012. The Biological Observation
421 Matrix (BIOM) format or: how I learned to stop worrying and love the ome-ome.
422 *GigaScience* 1:7.
- 423 McKinney W. 2012. pandas: a Foundational Python Library for Data Analysis and Statistics.
424 *O'Reilly Media, Inc*.
- 425 Mulkiyanian AY, Makarova KS, Galperin MY, and Koonin EV. 2007. Inventing the dynamo
426 machine: the evolution of the F-type and V-type ATPases. *Nat Rev Micro* 5:892-899.
- 427 Oinn T, Addis M, Ferris J, Marvin D, Senger M, Greenwood M, Carver T, Glover K, Pocock MR,
428 Wipat A, and Li P. 2004. Taverna: a tool for the composition and enactment of
429 bioinformatics workflows. *Bioinformatics* 20:3045-3054.
- 430 Ondov BD, Bergman NH, and Phillippy AM. 2011. Interactive metagenomic visualization in
431 a Web browser. *BMC bioinformatics* 12:385.

- 432 Qin J, Li R, Raes J, Arumugam M, Burgdorf K, Manichanh C, Nielsen T, Pons N, Levenez F,
433 Yamada T, Mende D, Li J, Xu J, Li S, Li D, Cao J, Wang B, Liang H, Zheng H, Xie Y, Tap J,
434 Lepage P, Bertalan M, Batto J, Hansen T, Le Paslier D, Linneberg A, Nielsen H,
435 Pelletier E, Renault P, Sicheritz-Ponten T, Turner K, Zhu H, Yu C, Li S, Jian M, Zhou Y,
436 Li Y, Zhang X, Li S, Qin N, Yang H, Wang J, Brunak S, Dore J, Guarner F, Kristiansen K,
437 Pedersen O, Parkhill J, Weissenbach J, Weissenbach J, Bork P, Ehrlich S, Wang J, and
438 Consortium M. 2010. A human gut microbial gene catalogue established by
439 metagenomic sequencing. *Nature* 464:59 - 65.
- 440 Ramirez KS, Leff JW, Barberán A, Bates ST, Betley J, Crowther TW, Kelly EF, Oldfield EE,
441 Shaw EA, and Steenbock C. 2014. Biogeographic patterns in below-ground diversity
442 in New York City's Central Park are similar to those observed globally. *Proceedings*
443 *of the Royal Society B: Biological Sciences* 281:20141988.
- 444 Razin S. 1992. Peculiar properties of mycoplasmas: The smallest self-replicating
445 prokaryotes. *FEMS Microbiology Letters* 100:423-431.
- 446 Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA,
447 Oakley BB, Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, Van Horn DJ, and
448 Weber CF. 2009. Introducing mothur: open-source, platform-independent,
449 community-supported software for describing and comparing microbial
450 communities. *Appl Environ Microbiol* 75:7537-7541.
- 451 Szczesnak A, Segata N, Qin X, Gevers D, Petrosino JF, Huttenhower C, Littman DR, and Ivanov,
452 II. 2011. The genome of th17 cell-inducing segmented filamentous bacteria reveals
453 extensive auxotrophy and adaptations to the intestinal environment. *Cell Host*
454 *Microbe* 10:260-272.
- 455 Segata N, Bornigen D, Morgan XC, and Huttenhower C. 2013. PhyloPhlAn is a new method
456 for improved phylogenetic and taxonomic placement of microbes. *Nat Commun*
457 4:2304.
- 458 Segata N, Haake SK, Mannon P, Lemon KP, Waldron L, Gevers D, Huttenhower C, and Izard J.
459 2012. Composition of the adult digestive tract bacterial microbiome based on seven
460 mouth surfaces, tonsils, throat and stool samples. *Genome Biol* 13:R42.
- 461 Segata N, Izard J, Waldron L, Gevers D, Miropolsky L, Garrett WS, and Huttenhower C. 2011.
462 Metagenomic biomarker discovery and explanation. *Genome Biol* 12:R60.
- 463 Shogan BD, Smith DP, Christley S, Gilbert JA, Zaborina O, and Alverdy JC. 2014. Intestinal
464 anastomotic injury alters spatially defined microbiome composition and function.
465 *Microbiome* 2:35.
- 466 Tree of Sex Consortium. 2014. Tree of Sex: A database of sexual systems. *Scientific Data* 1.
467 Wu GD, Chen J, Hoffmann C, Bittinger K, Chen YY, Keilbaugh SA, Bewtra M, Knights D,
468 Walters WA, Knight R, Sinha R, Gilroy E, Gupta K, Baldassano R, Nessel L, Li H,
469 Bushman FD, and Lewis JD. 2011. Linking long-term dietary patterns with gut
470 microbial enterotypes. *Science* 334:105-108.
- 471 Xu Z, Hansen MA, Hansen LH, Jacquiod S, and Sorensen SJ. 2014. Bioinformatic approaches
472 reveal metagenomic characterization of soil microbial community. *PLoS One*
473 9:e93445.

474 **Supplementary materials**
475

Figure 1 (on next page)

Figure 1. Schematic and simplified example of GraPhlAn visualization of annotated phylogenies and taxonomies.

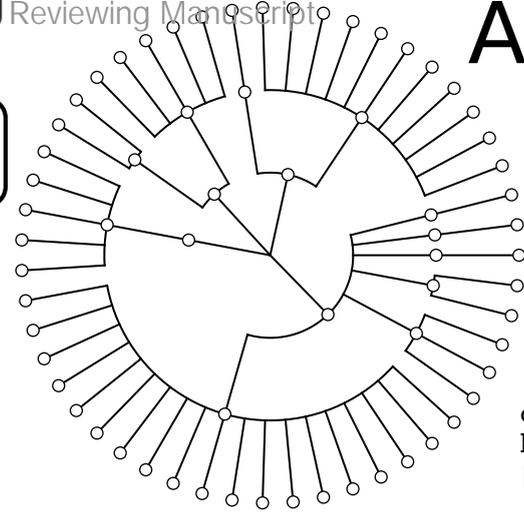
The software can start from a tree in Newick, Nexus, PhyloXML, or plain text formats. The “default plot” (A) produces a basic visualization of the tree's hierarchical structure. Through an annotation file, it is possible to configure a number of options that affect the appearance of the tree. For instance, some global parameters will affect the whole tree structure, such as the color and thickness of branches (“set global options”, B). The same annotation file can act on specific nodes, customizing their shape, size, and color (“set node options”, C). Labels and background colors for specific branches in the tree can also be configured (“set label options”, D). External to the circular area of the tree, the annotation file can include directives for plotting different shapes, heatmap colors, or bar-plots representing quantitative taxon traits (“set external ring options”, E).

Tree file
text, newick, or nexus formats

```
Bacillaceae.Bacillus.Bsubtilis
Bacillaceae.Bacillus.Bthuringiensis
[ . . . ]
Listeriaceae.Listeria.Lgrayi
Listeriaceae.Listeria.Linnocua
[ . . . ]
Paenibacillaceae.Brevibacillus.Bbrevis
Paenibacillaceae.Brevibacillus.Blaterosporus
[ . . . ]
Staphylococcaceae.Staphylococcus.Saureus
[ . . . ]
```

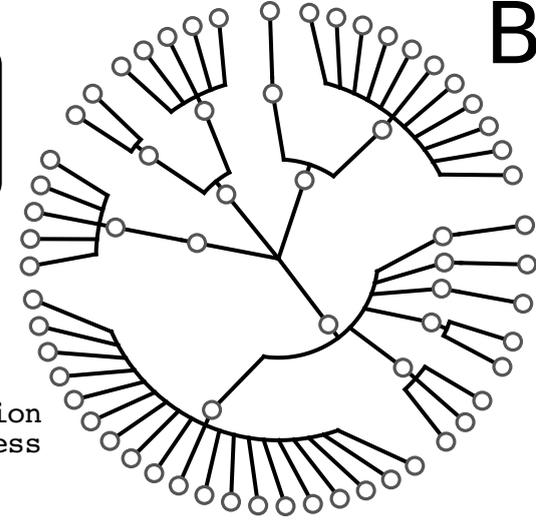
default plot

size
dpi
pad
format



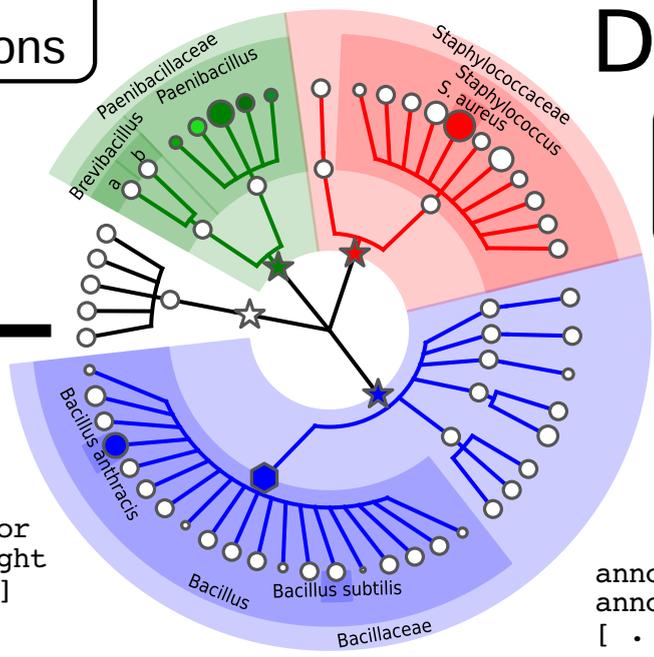
set global options

clade_separation
branch_thickness
[. . .]



set external ring options

ring_color
ring_height
[. . .]

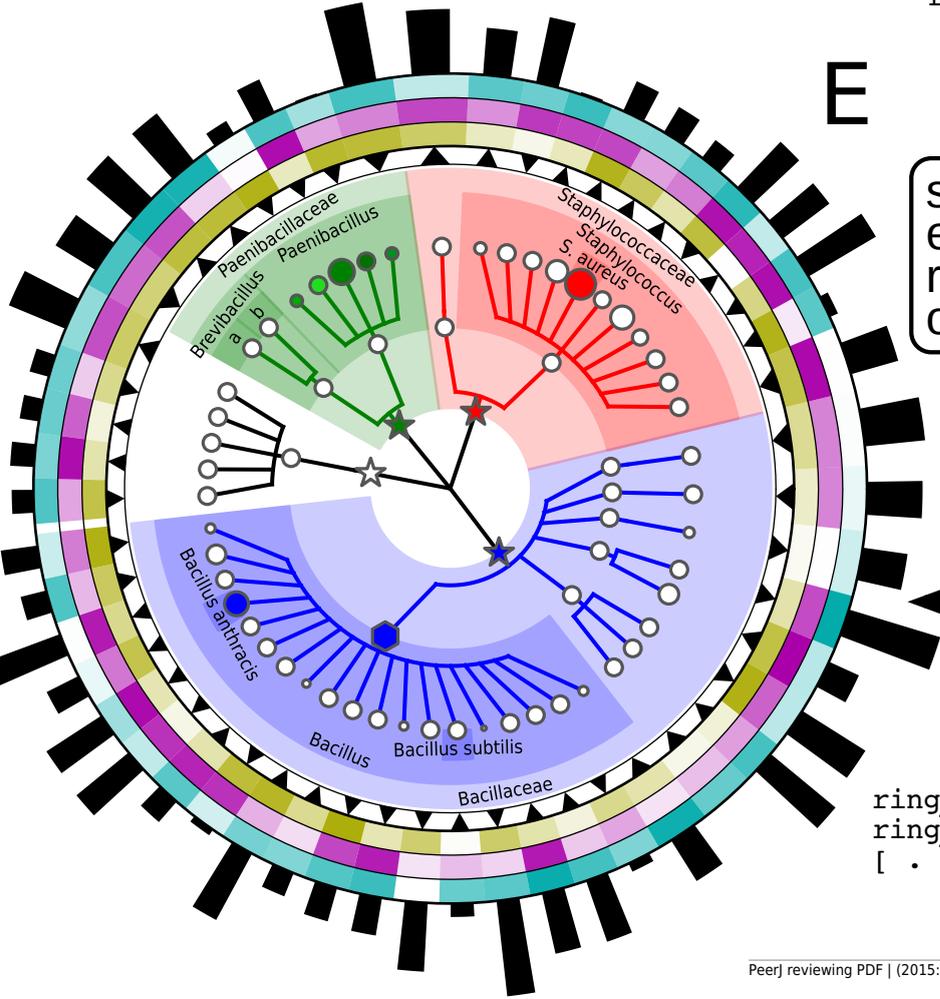
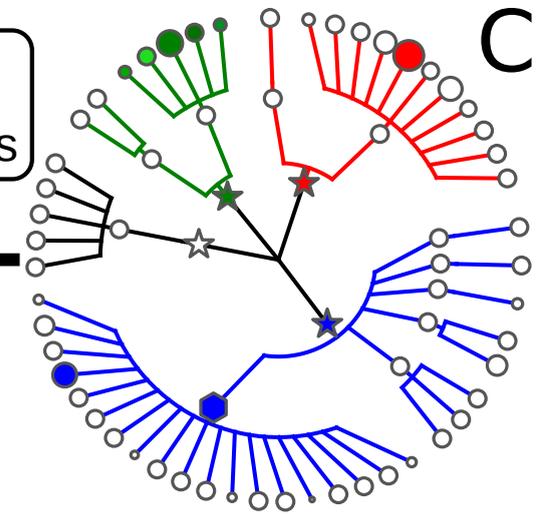


set node options

clade_marker_shape
clade_marker_size
[. . .]

set label options

annotation
annotation_background_color
[. . .]



2

Figure 2. A large, 3,737 genome phylogeny annotated with functional genomic properties.

We used the phylogenetic tree built using PhyloPhlAn (Segata et al. 2013) on all available microbial genomes as of 2013 and annotated the presence of ATP synthesis and Fatty Acid metabolism functional modules (as annotated in KEGG) and the genome length for all genomes. Colors and background annotation highlight bacterial phyla, and the functional information is reported in external rings. ATP synthesis rings visualize the presence (or absence) of each module, while Fatty Acid metabolism capability is represented with a gradient color. Data used in this image are available as indicated in the “Datasets used” paragraph, under “Materials and Methods” section.

A:Staphylococcus
 B:Enterococcus
 C:Streptococcus
 D:Lactobacillus
 E:Mycoplasma
 F:Helicobacter
 G:Campylobacter
 H:Burkholderia
 I:Xanthomonas
 J:Pseudomonas
 K:Shewanella
 L:Haemophilus
 M:Yersinia
 N:Escherichia
 O:Salmonella
 P:Enterobacter
 Q:Vibrio
 R:Corynebacterium
 S:Mycobacterium
 T:Bifidobacterium

ATP synthesis & Fatty Acid (FA) metabolism

● Actinobacteria
 ● Aquificae
 ● Bacteroidetes
 ● Chlamydiae
 ● Chlorobi
 ● Chloroflexi
 ● Crenarchaeota
 ● Cyanobacteria
 ● Euryarchaeota
 ● Firmicutes
 ● Proteobacteria
 ● Spirochaetes
 ● Tenericutes
 ● Thermi
 ● Thermotogae

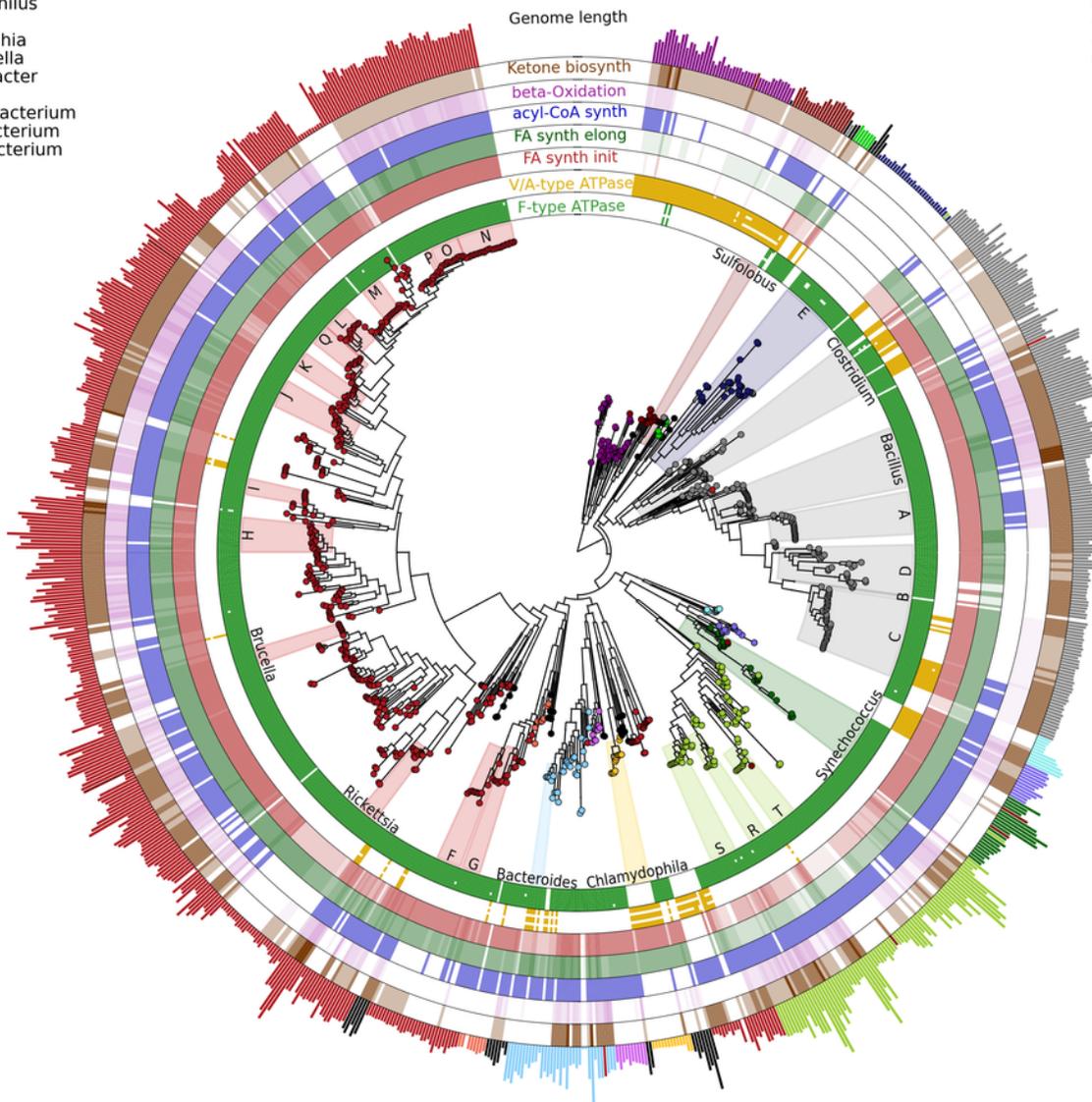


Figure 3 (on next page)

Figure 3. Taxonomic comparison between HMP and MetaHIT stool samples.

The taxonomic cladogram shows a comparison between the MetaHIT and HMP studies limited to samples from the gut (for the latter) and from healthy subjects (for the former). This image has been generated by GraPhlAn using input files from the supporting “export2graphlan” script (see “Materials and Methods”) applied on the output of MetaPhlAn2 (Segata et al. 2012b) and LEfSe (Segata et al. 2011) . Colors distinguish between HMP (green) and MetaHIT (blue), while the intensity reflects the LDA score, an indicator of the effect sizes of the significant differences. The size of the nodes correlates with their relative and logarithmically scaled abundances. Data used for this image is available as indicated under “Datasets used” paragraph in the “Materials and Methods” section.

MetaHIT vs. HMP (MetaPhlAn 2)

- A:Faecalibacterium
- B:Faecalibacterium prausnitzii
- C:Subdoligranulum
- D:Subdoligranulum unclassified
- E:Ruminococcus lactaris
- F:Ruminococcus sp 5 1 39BFAA
- G:Ruminococcus torques
- H:Coprococcus
- I:Coprococcus sp ART55 1
- J:Coprococcus comes
- K:Butyrivibrio
- L:Butyrivibrio crossotus
- M:Lachnospiraceae noname
- N:Dorea longicatena
- O:Roseburia hominis
- P:Roseburia inulinivorans
- Q:Eubacterium hallii
- R:Eubacterium siraeum
- S:Eubacterium eligens
- T:Clostridium sp L2 50
- U:Oscillibacter unclassified
- V:Erysipelotrichaceae
- W:Acidaminococcaceae
- X:Alistipes putredinis
- Y:Paraprevotella unclassified
- Z:Bacteroides coprocola
- a:Bacteroides caccae
- b:Bacteroides uniformis
- c:Bacteroides stercoris
- d:Bacteroides eggerthii
- e:Bacteroides sp 2 1 22
- f:Bacteroides ovatus
- g:Bacteroides thetaiotaomicron
- h:Bacteroides vulgatus
- i:Bacteroides xylanisolvens
- j:Bacteroides plebeius
- k:Barnesiella
- l:Barnesiella intestinihominis
- m:Parabacteroides unclassified
- n:Parabacteroides merdae
- o:Coriobacteriaceae
- p:Bifidobacteriaceae
- q:Bifidobacterium longum
- r:Bifidobacterium adolescentis

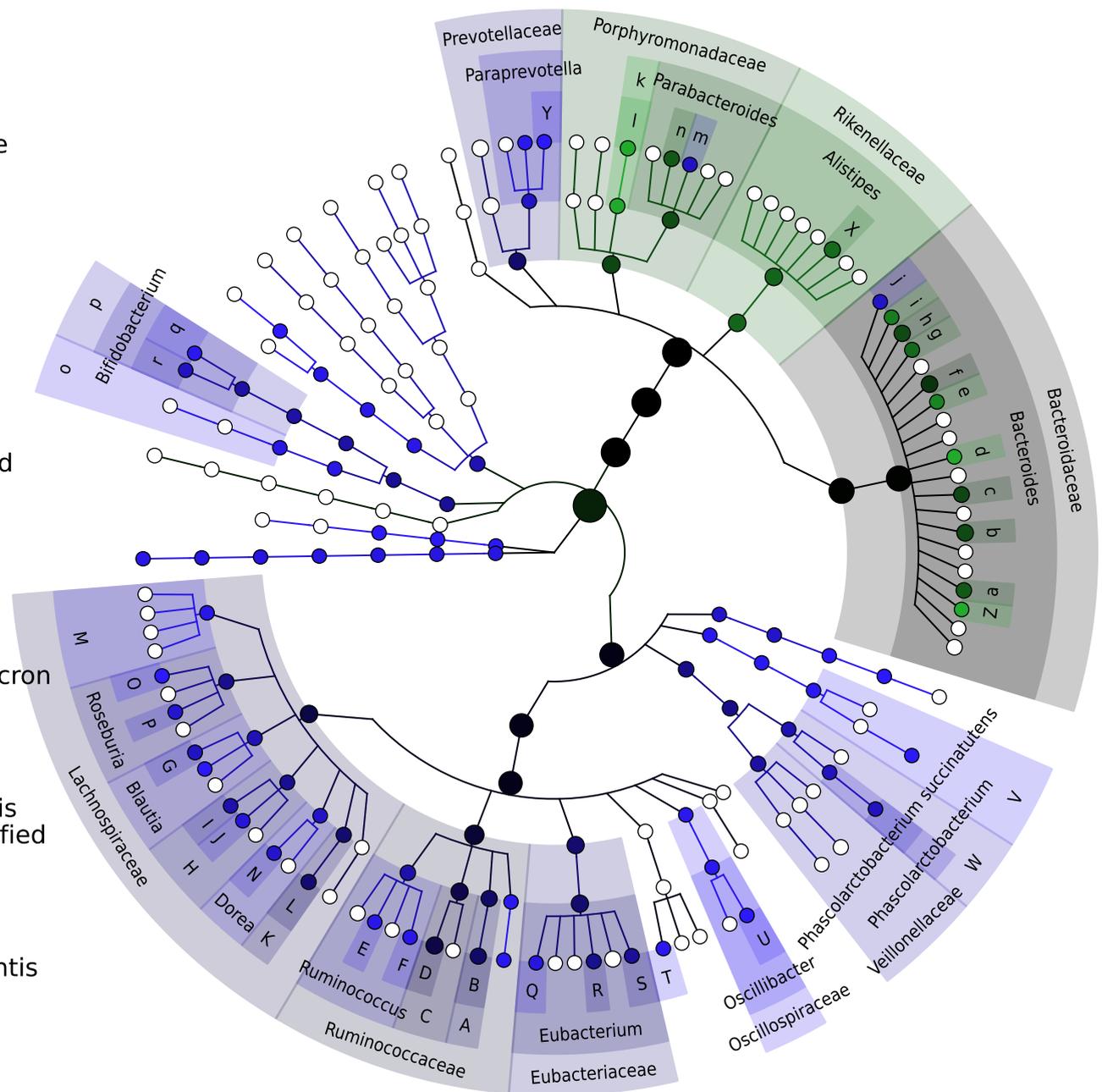


Figure 4(on next page)

Figure 4. Comparison of microbial community metabolic pathway abundances between HMP and MetaHIT.

*Comparison of functional pathway abundances from the HMP (green) and MetaHIT (blue). This is the functional counterpart of the plot in **Fig. 3** and was obtained applying GraPhlAn on HUMAnN (Abubucker et al. 2012) metabolic profiling. The intensity of the color represents the LDA score, and the sizes of the nodes are proportional to the pathway relative abundance estimated by HUMAnN. Three major groups are automatically highlighted by specifying them to the export2graphlan script: Environmental Information Processing, Genetic Information Processing, and Metabolism. Data used for this image is available as indicated under "Datasets used" paragraph in "Materials and Methods" section.*

Functional pathways

- A: Ribosome
- B: Sulfur relay system
- C: Protein export
- D: Homologous recombination
- E: Mismatch repair
- F: Base excision repair
- G: DNA replication
- H: Starch and sucrose metabolism
- I: Pentose and glucuronate interconversions
- J: Pentose phosphate pathway
- K: Glycolysis Gluconeogenesis
- L: Citrate cycle TCA cycle
- M: Drug metabolism other enzymes
- N: Pyrimidine metabolism
- O: Lipopolysaccharide biosynthesis
- P: Peptidoglycan biosynthesis
- Q: Lipoic acid metabolism
- R: Vitamin B6 metabolism
- S: Folate biosynthesis
- T: Biotin metabolism
- U: Nicotinate and nicotinamide metabolism
- V: Thiamine metabolism
- W: One carbon pool by folate
- X: Pantothenate and CoA biosynthesis
- Y: Terpenoid backbone biosynthesis
- Z: Streptomycin biosynthesis
- a: Sulfur metabolism
- b: Carbon fixation in photosynthetic organisms
- c: Cysteine and methionine metabolism
- d: Phenylalanine tyrosine and tryptophan biosynthesis
- e: Histidine metabolism
- f: Lysine biosynthesis
- g: Valine leucine and isoleucine biosynthesis
- h: Arginine and proline metabolism
- i: Alanine aspartate and glutamate metabolism
- j: Selenocompound metabolism
- k: D Glutamine and D glutamate metabolism
- l: D Alanine metabolism
- m: Bacterial secretion system
- n: Phosphotransferase system PTS

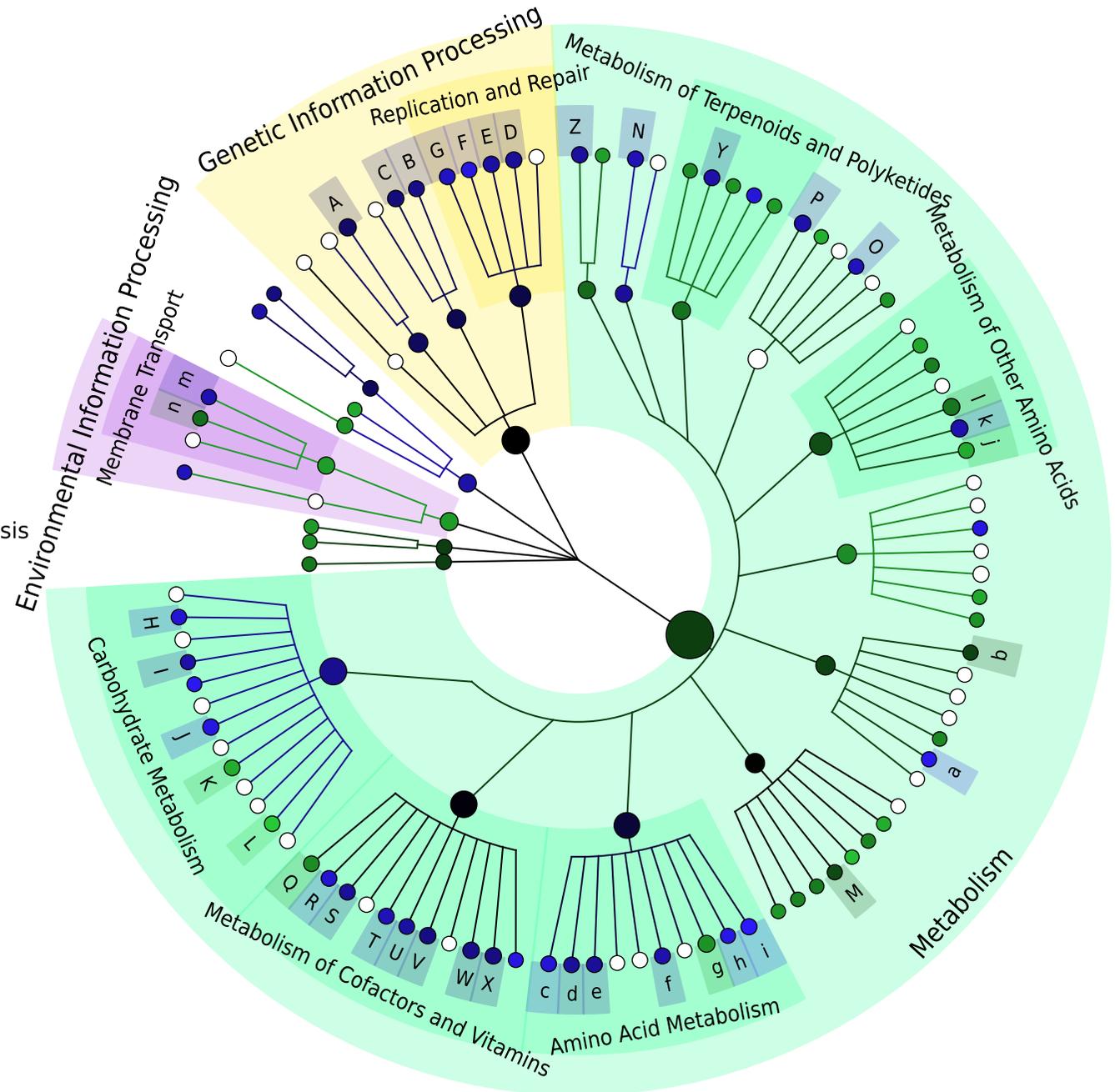


Figure 5 (on next page)

Figure 5. Integration of GraPhlAn into existing analyses pipelines.

We developed a conversion framework called “export2graphlan” that can deal with several output formats from different analysis pipelines, generating the necessary input files for GraPhlAn. Export2graphlan directly supports MetaPhlAn2, LEfSe, and HUMAnN output files. In addition, it can also accept BIOM files (both version 1 and 2), making GraPhlAn available for tools supporting this format including the QIIME and mothur systems. The tools can be ran on local machine as well as through the Galaxy web system using the modules reported in green boxes.

