

Manuscript Number:	GIGA-D-17-00230	
Full Title:	ASaiM: a Galaxy-based framework to analyze raw shotgun data from microbiota	
Article Type:	Technical Note	
Funding Information:	European Regional Development Fund	Not applicable
	Conseil Régional d'Auvergne	Dr Bérénice Batut
Abstract:	<p>New generation of sequencing platforms coupled to numerous bioinformatics tools has led to rapid technological progress in metagenomics and metatranscriptomics to investigate complex microorganism communities. Nevertheless, a combination of different bioinformatic tools remains necessary to draw conclusions out of microbiota studies. Modular and user-friendly tools would greatly improve such studies.</p> <p>We therefore developed ASaiM, an Open-Source Galaxy-based framework dedicated to microbiota data analyses. ASaiM provides a curated collection of tools to explore and visualize taxonomic and functional information from raw amplicon, metagenomic or metatranscriptomic sequences. To guide different analyses, several customizable workflows are included. All workflows are supported by tutorials and Galaxy interactive tours to guide the users through the analyses step by step. ASaiM is implemented as Galaxy Docker flavour. It is scalable to many thousand datasets, but also can be used a normal PC. The associated source code is available under Apache 2 license at https://github.com/ASaiM/framework and documentation can be found online (http://asaim.readthedocs.io/)</p> <p>Based on the Galaxy framework, ASaiM offers sophisticated analyses to scientists without command-line knowledge. ASaiM provides a powerful framework to easily and quickly explore microbiota data in a reproducible and transparent environment.</p>	
Corresponding Author:	Bérénice Batut, Ph.D. University of Freiburg Freiburg, GERMANY	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	University of Freiburg	
Corresponding Author's Secondary Institution:		
First Author:	Bérénice Batut, Ph.D.	
First Author Secondary Information:		
Order of Authors:	Bérénice Batut, Ph.D.	
	Kévin Gravouil	
	Clémence Defois	
	Saskia Hiltmann	
	Jean-François Brugère	
	Eric Peyretailade	
	Pierre Peyret	
Order of Authors Secondary Information:		
Opposed Reviewers:		
Additional Information:		

Question	Response
<p>Are you submitting this manuscript to a special series or article collection?</p>	<p>No</p>
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	<p>Yes</p>
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

1 ASaiM: a Galaxy-based framework to analyze raw shotgun
2 data from microbiota

3 B er enice Batut^{1,*}, K evin Gravouil^{2,3,4}, Cl emence Defois², Saskia Hiltermann⁵, Jean-Fran ois
4 Brug ere², Eric Peyretailade² and Pierre Peyret^{2,*}

5 Author affiliations

6 ¹Bioinformatics Group, Department of Computer Science, University of Freiburg, Germany

7 ²Universit e Clermont Auvergne, INRA, MEDIS, F-63000 Clermont-Ferrand, France

8 ³Universit e Clermont Auvergne, CNRS, LMGE, F-63000 Clermont-Ferrand, France

9 ⁴Universit e Clermont Auvergne, CNRS, LIMOS, F-63000 Clermont-Ferrand, France

10 ⁵Department of Bioinformatics, Erasmus University Medical Center, Rotterdam, 3015 CE,
11 Netherlands

12 Correspondence should be addressed to B.B. (berenice.batut@gmail.com) and P.P.
13 (pierre.peyret@uca.fr)

14 Abstract

15 Background

16 New generation of sequencing platforms coupled to numerous bioinformatics tools has led to
17 rapid technological progress in metagenomics and metatranscriptomics to investigate
18 complex microorganism communities. Nevertheless, a combination of different bioinformatic
19 tools remains necessary to draw conclusions out of microbiota studies. Modular and user-
20 friendly tools would greatly improve such studies.

21 Findings

22 We therefore developed ASaiM, an Open-Source Galaxy-based framework dedicated to
23 microbiota data analyses. ASaiM provides a curated collection of tools to explore and
24 visualize taxonomic and functional information from raw amplicon, metagenomic or
25 metatranscriptomic sequences. To guide different analyses, several customizable workflows
26 are included. All workflows are supported by tutorials and Galaxy interactive tours to guide
27 the users through the analyses step by step. ASaiM is implemented as Galaxy Docker
28 flavour. It is scalable to many thousand datasets, but also can be used a normal PC. The
29 associated source code is available under Apache 2 license at
30 <https://github.com/ASaiM/framework> and documentation can be found online
31 (<http://asaim.readthedocs.io>)

32 Conclusions

33 Based on the Galaxy framework, ASaiM offers sophisticated analyses to scientists without
34 command-line knowledge. ASaiM provides a powerful framework to easily and quickly
35 explore microbiota data in a reproducible and transparent environment.

36 Keywords

37 Metagenomics, Metabarcoding, User-friendly, Galaxy, Docker, Microbiota,

38 Findings

39 Background

40 The study of microbiota and microbial communities has been facilitated by the evolution of
41 sequencing techniques and the development of metagenomics and metatranscriptomics.
42 These techniques are giving insight into phylogenetic properties and metabolic components
43 of microbial communities. However, meta'omic data exploitation is not trivial due to the large
44 amount of data, high variability, incompleteness of reference databases, difficulty to find,
45 configure, use and combine the dedicated bioinformatics tools, etc. Hence, to extract useful
46 information, a sequenced microbiota sample has to be processed by sophisticated workflows
47 with numerous successive bioinformatics steps [1]. Each step may require execution of
48 several tools or software programs. For example, to extract taxonomic information with the
49 widely used QIIME [2] or Mothur [3], at least 10 different tools with at least 4 parameters
50 each are needed. Designed for amplicon data, both QIIME and Mothur can not be directly
51 applied to shotgun metagenomics data. In addition, the tools can be complex to use; they
52 are command-line tools and may require computational resources specially for the
53 metagenomics datasets. In this context, selecting the best tools, configuring them to use the
54 correct parameters and appropriate computational resources and combining them together
55 in an analysis chain is a complex and error-prone process. These issues and the involved
56 complexity are blocking scientist from participating in the analysis of their own data.
57 Furthermore, bioinformatics tools are often manually executed and/or patched together with
58 custom scripts. These practices raise doubts about a science gold standard: reproducibility
59 [3,4]. Web services and automated pipelines such as MG-RAST [5] and EBI metagenomics
60 [6] offer solutions to the accessibility issue. However, these web services work as a black

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84

box and are lacking in transparency, flexibility and even reproducibility as the version and parameters of the tools are not always available. Alternative approaches to improve accessibility, modularity and reproducibility can be found in open-source workflow systems such as Galaxy [6–8]. Galaxy is a lightweight environment providing a web-based, intuitive and accessible user interface to command-line tools, while automatically managing computation and transparently managing data provenance and workflow scheduling [6–8]. More than 4,500 tools can be used inside Galaxy environments. The tools can be selected and combined to build Galaxy flavors focusing on specific type of analysis, e.g. the Galaxy RNA workbench [9].

In this context, we developed ASaiM (Auvergne Sequence analysis of intestinal Microbiota), an Open-Source opinionated Galaxy-based framework. It integrates tools and workflows dedicated to microbiota analyses with an extensive documentation (<http://asaim.readthedocs.org>) and training support.

Goals of ASaiM

ASaiM is developed as a modular, accessible, redistributable, sharable and user-friendly framework for scientists working with microbiota data. This framework is unique in combining curated tools and workflows and providing easy access for scientists.

ASaiM is based on four pillars: 1) easy and stable dissemination via Galaxy, Docker and conda, 2) a comprehensive set of metagenomics related tools, 3) a set of predefined and tested workflows, and 4) extensive documentation and training to help scientists in their analyses.

A framework built on the shoulders of giants

The ASaiM framework is built on existing tools and infrastructures and combine all their forces to build an easily accessible and reproducible analysis platform.

1 85 ASaiM is implemented as portable virtualized container based on Galaxy framework [8].
2 86 Galaxy provides researchers with means to reproduce their own workflows analyses, rerun
3
4 87 entire pipelines, or publish and share them with others. Based on Galaxy, ASaiM is scalable
5
6 88 from single CPU installations to large multi-node high performance computing environments.
7
8 89 Deployments can be archived by using a pre-built ASaiM Docker image, which is based on
9
10 90 the Galaxy Docker project (<http://bgruening.github.io/docker-galaxy-stable>) or by installing all
11
12 91 needed components into an already existing Galaxy instance. This ASaiM Docker instance
13
14 92 is customized with a variety of selected tools, workflows, Interactive tours and data that have
15
16 93 been added as additional layers on top of the generic Galaxy Docker instance. The
17
18 94 containerization keeps the deployment task to a minimum. The selected Galaxy tools are
19
20 95 automatically installed from the Galaxy ToolShed [10] (<https://toolshed.g2.bx.psu.edu/>) using
21
22 96 the Galaxy API BioBlend [11] and the installation of the tools and their dependencies are
23
24 97 automatically resolved using packages available through Bioconda
25
26 98 (<https://bioconda.github.io>). We migrated then 10 tools/suites of tools and their
27
28 99 dependencies to Bioconda (e.g. HUMAnN2) and integrated 14 suites into Galaxy (e.g.
29
30 100 QIIME with around forty tools).
31
32
33 101 The containerization as well as the packaging with conda enables automatic continuous
34
35 102 integration tests at different levels: dependencies (BioConda), tool integration in Galaxy,
36
37 103 Galaxy itself and at ASaiM level. Together with strict version management on all levels, this
38
39 104 contributes to a high degree of error-control and reproducibility.
40
41
42
43
44
45
46
47

48 106 Tools for microbiota data analyses

49
50
51 107 The tools integrated in ASaiM can be seen in Table 1. They are expertly selected for their
52
53 108 relevance with regard to microbiota studies, such as Mothur [3], QIIME [2], MetaPhlan2 [12],
54
55 109 HUMAnN2 [13] or tools used in existing pipelines such as EBI Metagenomics' one. We also
56
57 110 added general tools used in sequence analysis such as quality control, mapping or similarity
58
59 111 search tools.
60
61
62
63
64
65

112 **Table 1:** Available tools in ASaiM

Section	Subsection	Tools
File and meta tools	Data retrieval	EBISearch, ENASearch, SRA Tools
	Text manipulation	Tools from Galaxy ToolShed
	Sequence file manipulation	Tools from Galaxy ToolShed
	BAM/SAM file manipulation	SAM tools [14–16]
	BIOM file manipulation	BIOM-Format tools [17]
Genomics tools	Assembly	FastQ joiner [18], FastQ-join
	Quality control	FastQC , PRINSEQ [19], Trim Galore! , Trimmomatic [20], MultiQC [21]
	Clustering	CD-Hit [22], Format CD-HIT outputs
	Sorting and prediction	SortMeRNA [23], FragGeneScan [24]
	Mapping	BWA [25,26], Bowtie [27]
	Similarity search	NCBI Blast+ [28,29], Diamond [30]
	Alignment	HMMER3
Microbiota dedicated tools	Metagenomics data manipulation	VSEARCH [31]
	Amplicon sequence processing	Mothur [3], QIIME [2]
	Taxonomy assignment on WGS sequences	MetaPhlan2 [12], Format MetaPhlan2, Kraken [32]
	Metabolism assignment	HUMAN2 [13], Group HUMAN2 to GO slim terms , Compare HUMAN2 outputs, PICRUST [33], InterProScan
	Combination of functional and taxonomic results	Combine MetaPhlan2 and HUMAN2 outputs
	Visualization	Export2graphlan, GraPhlan [34], KRONA [35]

113 This table presents the tools, organized in section and subsections to help users. A more detailed
 114 table of the available tools and some documentation can be found in the online documentation
 115 (<http://asaim.readthedocs.io>)

116

1
2 117 An effort in development was made to integrate these tools into Conda and the Galaxy
3
4 118 environment, with the help and support of the Galaxy community. We also developed two
5
6 119 new tools to search and get data from EBI Metagenomics and ENA databases using the API
7
8
9 120 of the databases (EBISearch and ENASearch) and a tool to group HUMAnN2 outputs into
10
11 121 Gene Ontology Slim Terms. Tools inside ASaiM are organized to make them findable and
12
13 122 documented (<http://asaim.readthedocs.io>).

16 123 Diverse source of data

17
18
19 124 Any easy way to upload user-data into ASaiM is provided by an web-interface or more
20
21 125 sophisticated via FTP or SFTP. Moreover, we added specialised tools that can interact with
22
23 126 external databases like NCBI, ENA or EBI Metagenomics to query them and download data
24
25
26 127 into the framework.

29 128 Visualization of the data

30
31
32 129 An analysis often ends with summarizing figures that conclude and represent the findings.
33
34 130 ASaiM includes standard interactive plotting tools to draw bar charts and scatter plots from
35
36
37 131 all kinds of tabular data. Phinch visualization is also included to interactively visualize and
38
39 132 explore any BIOM file, and generate different types of ready-to-publish figures. We also
40
41 133 integrated two other tools to explore and represent the community structure from outputs of
42
43 134 MetaPhlAn: KRONA [35] and GraPhlAn. Moreover, as in any Galaxy instance, other
44
45 135 visualization are included such Phyloviz for phylogenetic trees or the Genome browser
46
47
48 136 Trackster for visualizing SAM/BAM, BED, GFF/GTF, WIG, bigWig, bigBed, bedGraph, and
49
50 137 VCF datasets.

54 138 Workflows

55
56
57 139 Each tool can be used separately in an explorative manner or multiple tools can be
58
59 140 orchestrated inside workflows passing raw data to the data reduction step, to information

141 extraction and visualization. To assist in microbiota analyses, several default but
142 customizable workflows are proposed in ASaiM. All the available workflows with tool and
143 parameter choices are documented (<http://asaim.readthedocs.io>).

144 Analysis of raw metagenomic or metatranscriptomic shotgun data

145 A workflow quickly produces, from raw metagenomic or metatranscriptomic shotgun data,
146 accurate and precise taxonomic assignments, wide extended functional results and
147 taxonomically related metabolism information (Figure 1). This workflow consists of i)
148 processing with quality control/trimming (FastQC and Trim Galore!) and dereplication
149 (VSearch [31]; ii) taxonomic analyses with assignment (MetaPhlAn2 [12]) and visualization
150 (KRONA , GraPhlAn); iii) functional analyses with metabolic assignment and pathway
151 reconstruction (HUMAN2 [13]); iv) functional and taxonomic combination with developed
152 tools combining HUMAN2 and MetaPhlAn2 outputs.

153 This workflow has been tested on two mock metagenomic datasets with controlled
154 communities (Supplementary material). We have compared the extracted taxonomic and
155 functional information to such information extracted with the EBI metagenomics' pipeline and
156 to the expectations from the mock datasets. With ASaiM, we generate more accurate and
157 precise data for taxonomic analyses (Figure 2): we can access information at the level of the
158 species. More informative data for metabolic description (gene families, gene ontologies,
159 pathways, etc) are also extracted with ASaiM compared to the ones available on EBI
160 metagenomics. With this workflow, we can investigate which taxons are involved in a
161 specific pathway or a gene family (e.g. involved species and their relative involvement in
162 different step of fatty acid biosynthesis pathways, Figure 3).

163 For the tests, ASaiM was deployed on a computer with Debian GNU/Linux System, 8 cores
164 Intel(R) Xeon(R) at 2.40 GHz and 32 Go of RAM. The workflow processed the 1,225,169
165 and 1,386,198 454 GS FLX Titanium reads of each datasets in 4h44 and 5h22 respectively,
166 with a stable memory usage (Supplementary material). With this workflow, it is then easy
167 and quick to process raw microbiota data and extract diverse useful information.

168 Analysis of amplicon data

1
2
3 169 To analyze amplicon data, the Mothur and QIIME tool suites are available to ASaiM. We
4
5 170 integrated the workflows described in tutorials of Mothur and QIIME websites, as example of
6
7 171 amplicon data analyses as well as support for the training material. These workflows, as any
8
9 172 workflows available in ASaiM, can be adapted for a specific analysis or used as
10
11 173 subworkflows by the users.

15 174 Running as in EBI metagenomics

16
17
18 175 The tools used in the EBI Metagenomics pipeline are also available in ASaiM. We integrate
19
20 176 then also a workflow with the same steps as the EBI Metagenomics pipeline. Analyses made
21
22 177 in EBI Metagenomics website can be then reproduced locally, without having to wait for
23
24 178 availability of EBI Metagenomics or to upload any data on EBI Metagenomics. However the
25
26 179 parameters must be defined by the user as we can not find them on EBI Metagenomics
27
28
29 180 documentation.

33 181 Documentation and training

34
35
36 182 A tool or software is easier to use if it is well documented. Hence extensive documentation
37
38 183 helps the users to be familiar with the tool and also prevents mis-usage. For ASaiM, we
39
40 184 developed an extensive online documentation (<http://asaim.readthedocs.io>), mainly to
41
42 185 explain how to use it, how to deploy it, which tools are integrated with small documentation
43
44 186 about these tools, which workflows are integrated and how to use them.

45
46
47 187 In addition to this online documentation, Galaxy Interactive Tours are included inside the
48
49 188 Galaxy instance. Such tours guide users through an entire analysis in an interactive (step-
50
51 189 by-step) way. Some tours, included in every Galaxy instance, explains how to use Galaxy.
52
53 190 We also developed such tours dedicated specifically to the ASaiM workflows.

54
55
56 191 These interactive tours are used to complement tutorials and trainings. Some tutorials about
57
58 192 the integrated workflows have been developed to explain step-by-step the workflows with
59
60
61
62
63
64
65

193 small example datasets. Hosted in the Galaxy Training Network (GTN) GitHub repository
194 (<https://github.com/galaxyproject/training-material>), the tutorials are available online at
195 <http://training.galaxyproject.org>. They have been used during several workshops on
196 metagenomics data analysis with ASaiM as training support. These tutorials are also
197 accessible directly from ASaiM and its documentation for self-training.

198 Installation and running ASaiM

199 Running the containerized ASaiM simply requires to install Docker and to start the ASaiM
200 image with:

```
201 $ docker run -d -p 8080:80 quay.io/bebatut/asaim
```

202 Thanks to Docker, ASaiM can be installed under every operating systems, even with a
203 graphical tool (Kitematic: <https://kitematic.com>) on OSX and Windows.

204 ASaiM is production-ready. It can also be configured to use external accessible computer
205 clusters or cloud environments.

206 It is also possible and easy to install all or only a subset of tools of the ASaiM framework on
207 existing Galaxy instances. The set of available tools can be easily extended either only a
208 given instance using the Galaxy admin interface or for ASaiM more generally thanks to the
209 simple definition of the installed tools in YAML files available in ASaiM GitHub repository. In
210 the latter case, the Docker image will be automatically rebuilt and the already integrated
211 tools will be updated to keep ASaiM up-to-date. For reproducibility reason, every version of
212 the Docker image is associated to a tag and is conserved.

213 Conclusion

214 ASaiM provides a powerful framework to easily and quickly analyze microbiota data in a
215 reproducible, accessible and transparent way. Built on a Galaxy instance wrapped in a
216 Docker image, ASaiM can be easily deployed with a comprehensive set of tools and their
217 dependencies. These tools are complemented with a set of predefined and tested workflows

218 to address the main microbiota questions (community structure and functions). All these
219 tools and workflows are extensively documented online (<http://asaim.readthedocs.io>) and
220 supported by Galaxy Interactive Tours and tutorials.

221 With this complete infrastructure, ASaiM offers a good environment for sophisticated
222 microbiota analyses to scientists without computational knowledge, while promoting
223 transparency, sharing and reproducibility.

224 Methods

225 For the tests, ASaiM was deployed on a computer with Debian GNU/Linux System, 8 cores
226 Intel(R) Xeon(R) at 2.40 GHz and 32 Go of RAM. The workflow has been run on two mock
227 community samples of Human Microbiome Project (HMP), containing a genomic mixture of
228 22 known microbial strains. The details of comparison analyses are described in the
229 Supplementary Material.

230 Availability of supporting source code and requirements

- 231 ● Project name: ASaiM
- 232 ● Project home page: <https://github.com/ASaiM/framework>
- 233 ● Operating system(s): Platform independent
- 234 ● Other requirements: Docker
- 235 ● License: Apache 2

236 All tools described herein are available in the Galaxy Toolshed
237 (<https://toolshed.g2.bx.psu.edu>). The Dockerfile to automatically install deploy ASaiM is
238 provided in the GitHub repository and a pre-built Docker image is available at
239 <https://quay.io/repository/bebatut/asaim-framework>.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

240 Declarations

241 Competing interests

242 The author(s) declare that they have no competing interests.

243 Funding

244 The Auvergne Regional Council and the European Regional Development Fund have
245 supported this work.

246 Authors' contributions

247 BB, KG, CD, SH, JFB, EP, PP contributed equally to the conceptualization, to the
248 methodology and to the writing process. JFP, PP contributed equally to the funding
249 acquisition. BB, KG, SH contributed equally to the software development and BB, KG, CD
250 and JFP to the validation.

251 Acknowledgements

252 The authors would like to thank EA 4678 CIDAM, UR 454 INRA, M2iSH, LIMOS, CRR1,
253 de.NBI for their involvement in this project, as well as Réjane Beugnot, Thomas Eymard,
254 David Parsons and Björn Grüning for their help.

255 References

- 256 1. Ladoukakis E, Kolisis FN, Chatziioannou AA. Integrative workflows for metagenomic
257 analysis. *Front Cell Dev Biol.* 2014;2:70.
- 258 2. Kuczynski J, Stombaugh J, Walters WA, González A, Gregory Caporaso J, Knight R.
259 Using QIIME to Analyze 16S rRNA Gene Sequences from Microbial Communities. *Current*
260 *Protocols in Microbiology.* 2012. p. 1E.5.1–1E.5.20.
- 261 3. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, et al. Introducing
262 mothur: open-source, platform-independent, community-supported software for describing
263 and comparing microbial communities. *Appl. Environ. Microbiol.* 2009;75:7537–41.

264 4. Nekrutenko A, Taylor J. Next-generation sequencing data interpretation: enhancing
1 265 reproducibility and accessibility. *Nat. Rev. Genet.* 2012;13:667–72.
2
3 266 5. Meyer F, Paarmann D, D'Souza M, Olson R, Glass EM, Kubal M, et al. The
4 267 metagenomics RAST server – a public resource for the automatic phylogenetic and
5 268 functional analysis of metagenomes. *BMC Bioinformatics.* 2008;9:386.
6
7 269 6. Hunter S, Corbett M, Denise H, Fraser M, Gonzalez-Beltran A, Hunter C, et al. EBI
8 270 metagenomics--a new resource for the analysis and archiving of metagenomic data. *Nucleic
9 271 Acids Res.* 2014;42:D600–6.
10
11 272 7. Goecks J, Nekrutenko A, Taylor J, Galaxy Team. Galaxy: a comprehensive approach for
12 273 supporting accessible, reproducible, and transparent computational research in the life
14 274 sciences. *Genome Biol.* 2010;11:R86.
15
16 275 8. Afgan E, Baker D, van den Beek M, Blankenberg D, Bouvier D, Čech M, et al. The Galaxy
17 276 platform for accessible, reproducible and collaborative biomedical analyses: 2016 update.
18 277 *Nucleic Acids Res.* 2016;44:W3–10.
19
20 278 9. Grüning BA, Fallmann J, Yusuf D, Will S, Erxleben A, Eggenhofer F, et al. The RNA
21 279 workbench: best practices for RNA and high-throughput sequencing bioinformatics in
22 280 Galaxy. *Nucleic Acids Res.* [Internet]. 2017; Available from:
23 281 <http://dx.doi.org/10.1093/nar/gkx409>
24
25 282 10. Blankenberg D, Von Kuster G, Bouvier E, Baker D, Afgan E, Stoler N, et al.
27 283 Dissemination of scientific software with Galaxy ToolShed. *Genome Biol.* 2014;15:403.
28
29 284 11. Sloggett C, Goonasekera N, Afgan E. BioBlend: automating pipeline analyses within
30 285 Galaxy and CloudMan. *Bioinformatics.* 2013;29:1685–6.
31
32 286 12. Truong DT, Franzosa EA, Tickle TL, Scholz M, Weingart G, Pasolli E, et al. MetaPhlan2
33 287 for enhanced metagenomic taxonomic profiling. *Nat. Methods.* 2015;12:902–3.
34
35 288 13. Abubucker S, Segata N, Goll J, Schubert AM, Izard J, Cantarel BL, et al. Metabolic
36 289 reconstruction for metagenomic data and its application to the human microbiome. *PLoS
38 290 Comput. Biol.* 2012;8:e1002358.
39
40 291 14. Li H. A statistical framework for SNP calling, mutation discovery, association mapping
41 292 and population genetical parameter estimation from sequencing data. *Bioinformatics.*
42 293 2011;27:2987–93.
43
44 294 15. Li H. Improving SNP discovery by base alignment quality. *Bioinformatics.* 2011;27:1157–
45 295 8.
46
47 296 16. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
48 297 Alignment/Map format and SAMtools. *Bioinformatics.* 2009;25:2078–9.
49
50 298 17. McDonald D, Clemente JC, Kuczynski J, Rideout JR, Stombaugh J, Wendel D, et al. The
51 299 Biological Observation Matrix (BIOM) format or: how I learned to stop worrying and love the
52 300 ome-ome. *Gigascience.* 2012;1:7.
53
54 301 18. Blankenberg D, Gordon A, Von Kuster G, Coraor N, Taylor J, Nekrutenko A, et al.
55 302 Manipulation of FASTQ data with Galaxy. *Bioinformatics.* 2010;26:1783–5.
56
57 303 19. Schmieder R, Edwards R. Quality control and preprocessing of metagenomic datasets.
58 304 *Bioinformatics.* 2011;27:863–4.
59
60
61
62
63
64
65

305 20. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence
1 306 data. *Bioinformatics*. 2014;30:2114–20.
2

3 307 21. Ewels P, Magnusson M, Lundin S, Käller M. MultiQC: summarize analysis results for
4 308 multiple tools and samples in a single report. *Bioinformatics*. 2016;32:3047–8.
5

6 309 22. Fu L, Niu B, Zhu Z, Wu S, Li W. CD-HIT: accelerated for clustering the next-generation
7 310 sequencing data. *Bioinformatics*. 2012;28:3150–2.
8

9 311 23. Kopylova E, Noé L, Touzet H. SortMeRNA: fast and accurate filtering of ribosomal RNAs
10 312 in metatranscriptomic data. *Bioinformatics*. 2012;28:3211–7.
11

12 313 24. Rho M, Tang H, Ye Y. FragGeneScan: predicting genes in short and error-prone reads.
13 314 *Nucleic Acids Res*. 2010;38:e191–e191.
14

15 315 25. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.
16 316 *Bioinformatics*. 2009;25:1754–60.
17

18 317 26. Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler transform.
19 318 *Bioinformatics*. 2010;26:589–95.
20

21 319 27. Press AR. Ultrafast and Memory-Efficient Alignment of Short DNA Sequences to the
22 320 Human Genome. CreateSpace; 2015.
23

24 321 28. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:
25 322 architecture and applications. *BMC Bioinformatics*. 2009;10:421.
26

27 323 29. Cock PJA, Chilton JM, Grüning B, Johnson JE, Soranzo N. NCBI BLAST+ integrated
28 324 into Galaxy. *Gigascience*. 2015;4:39.
29

30 325 30. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND.
31 326 *Nat. Methods*. 2015;12:59–60.
32

33 327 31. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. VSEARCH: a versatile open source
34 328 tool for metagenomics. *PeerJ*. 2016;4:e2584.
35

36 329 32. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using
37 330 exact alignments. *Genome Biol*. 2014;15:R46.
38

39 331 33. Langille MGI, Zaneveld J, Caporaso JG, McDonald D, Knights D, Reyes JA, et al.
40 332 Predictive functional profiling of microbial communities using 16S rRNA marker gene
41 333 sequences. *Nat. Biotechnol*. 2013;31:814–21.
42

43 334 34. Asnicar F, Weingart G, Tickle TL, Huttenhower C, Segata N. Compact graphical
44 335 representation of phylogenetic data and metadata with GraPhlAn. *PeerJ*. 2015;3:e1029.
45

46 336 35. Ondov BD, Bergman NH, Phillippy AM. Interactive metagenomic visualization in a Web
47 337 browser. *BMC Bioinformatics*. 2011;12:385.
48

49 338
50
51
52
53
54 339
55
56
57
58
59
60
61
62
63
64
65

340 **Figure 1:** Main ASaiM workflow to analyze raw sequences.

341 This workflow takes as input a dataset of raw shotgun sequences (in FastQ format) from
342 microbiota, preprocess it (yellow boxes), extracts taxonomic (red boxes) and functional
343 (purple boxes) assignments and combines them (green boxes).

344 Image available under CC-BY license (<https://doi.org/10.6084/m9.figshare.5371396.v3>)

345

346 **Figure 2:** Comparisons of the community structure for SRR072233.

347 This figure compares the community structure between the expectations (mapping of the
348 sequences on the expected genomes), data found on EBI Metagenomics database
349 (extracted with the EBI Metagenomics pipeline) and the results of the main ASaiM workflow
350 (Figure 1).

351

352 **Figure 3:** Example of an investigation of the relation between community structure and
353 functions.

354 The involved species and their relative involvement in fatty acid biosynthesis pathways have
355 been extracted with ASaiM workflow (Figure 1) for SRR072233

356

TAXONOMIC ANALYSES

PROCESSING

FUNCTIONAL ANALYSES

FUNCTIONAL AND TAXONOMIC COMBINATION

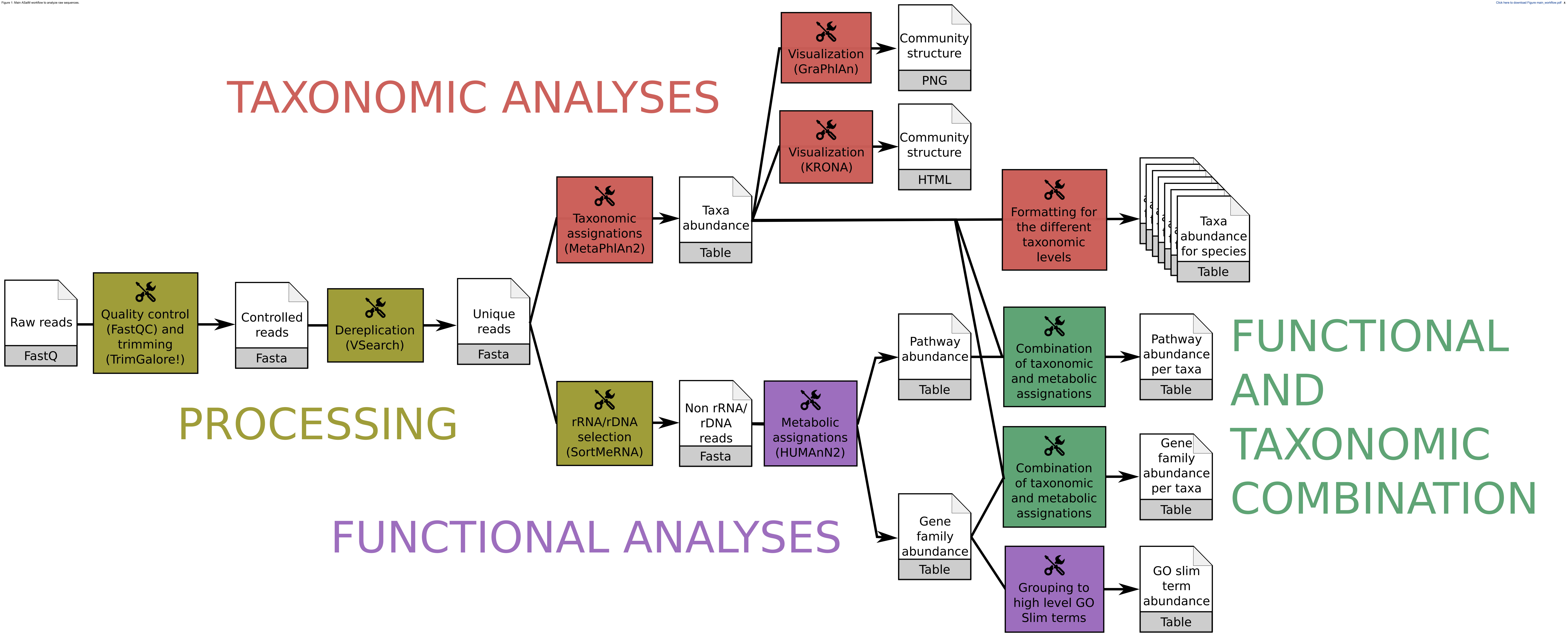


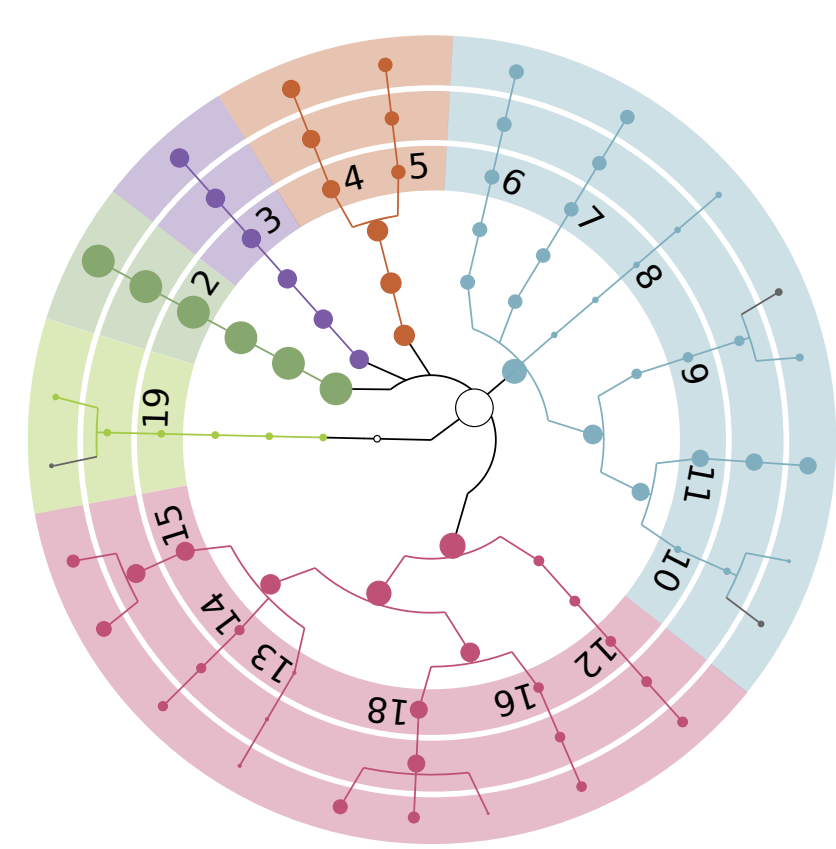
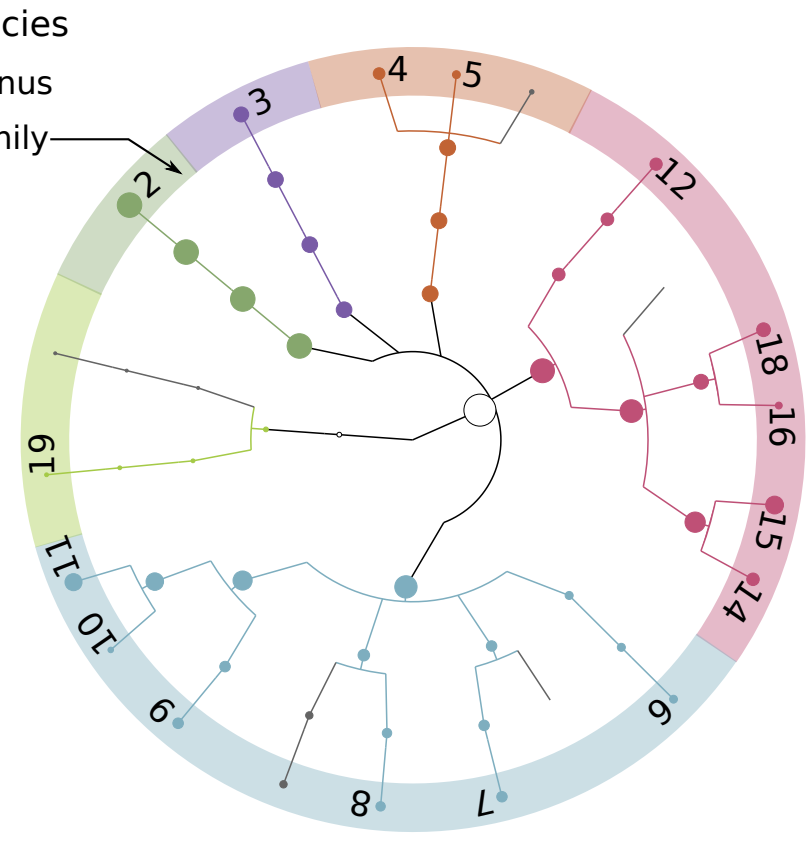
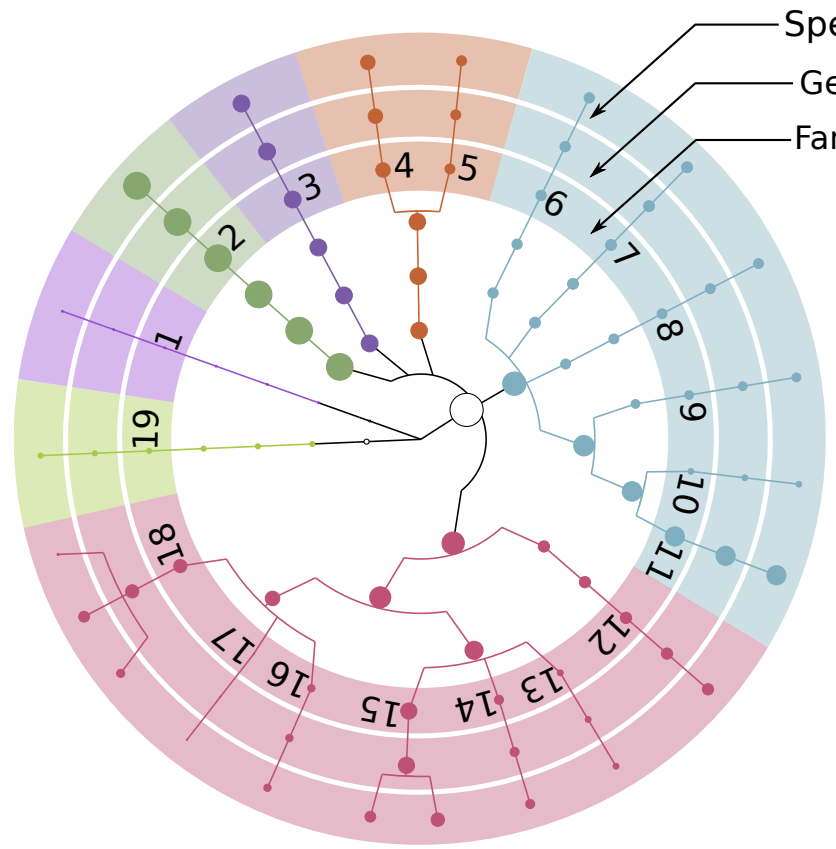
Figure 2: Comparisons of the community structure for SRR072233

[Click here to download Figure hmp_taxonomic_results.pdf](#)

Expectations

EBI Metagenomics results

ASaiM framework results



Phyla

Ascomycota	Proteobacteria
Deinococcus-Thermus	Firmicutes
Bacteroidetes	Euryarchaeota
Actinobacteria	Unexpected

Families

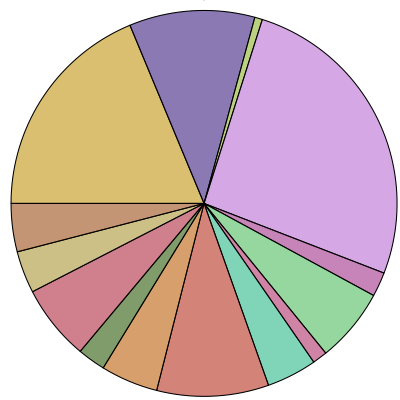
1 Debaryomycetaceae	5 Actinomycetaceae	9 Enterobacteriaceae	13 Bacillaceae	17 Lactobacillaceae
2 Deinococcaceae	6 Helicobacteraceae	10 Pseudomonadaceae	14 Listeriaceae	18 Streptococcaceae
3 Bacteroidaceae	7 Neisseriaceae	11 Moraxellaceae	15 Staphylococcaceae	19 Methanobacteriaceae
4 Propionibacteriaceae	8 Rhodobacteraceae	12 Clostridiales	16 Enterococcaceae	

Figure 3: Example of an investigation of the relationship between community structure and functions

[Click here to download Figure hmp_taxonomically_related_functional_results.pdf](#)

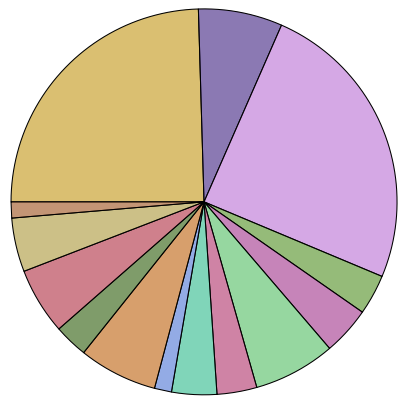


Superpathway of fatty acid biosynthesis initiation (FASYN-INITIAL-PWY)



an acetoacetyl-acp

Pathway of fatty acid elongation (FASYN-ELONG-PWY)



Species

- Acinetobacter baumannii
- Bacteroides vulgatus
- Clostridium beijerinckii
- Deinococcus radiodurans
- Enterococcus faecalis
- Escherichia coli
- Helicobacter pylori
- Listeria monocytogenes
- Neisseria meningitidis
- Propionibacterium acnes
- Pseudomonas aeruginosa
- Rhodobacter sphaeroides
- Staphylococcus aureus
- Staphylococcus epidermidis
- Streptococcus mitis oralis pneumoniae
- Streptococcus mutans



Click here to access/download
Supplementary Material
report.pdf

