
MG7: Configurable and scalable 16S metagenomics data analysis – new methods optimized for massive cloud computing

Alexey Alekhin^{† 1} Evdokim Kovach^{† 1} Marina Manrique¹ Pablo Pareja-Tobes¹
Eduardo Pareja¹ Raquel Tobes¹ and Eduardo Pareja-Tobes^{1,*}

¹*Oh no sequences! Research Group, Era7 Bioinformatics, Granada, Spain*

Correspondence*:

Eduardo Pareja-Tobes

Oh no sequences! Research Group, Era7 Bioinformatics, Plaza Campo Verde 3,
Granada, 18001, Spain, eparejatobes@ohnosequences.com

2 ABSTRACT

3 As part of the Cambrian explosion of omics data, metagenomics brings to the table a specific,
4 defining trait: its social essence. The *meta* prefix exerts its influence, with multitudes manifesting
5 themselves everywhere; from samples to data analysis, from actors involved to (present and
6 future) applications. Of these dimensions, data analysis is where needs lay further from what
7 current tools provide. Key features are, among others, scalability, reproducibility, data provenance
8 and distribution, process identity and versioning. These are the goals guiding our work in
9 MG7, a 16S metagenomics data analysis system. The basic principle is a new approach to
10 data analysis, where configuration, processes, or data locations are static, type-checked and
11 subject to the standard evolution of a well-maintained software project. Cloud computing, in its
12 Amazon Web Services incarnation, when coupled with these ideas, produces a robust, safely
13 configurable, scalable tool. Processes, data, machine behaviors and their dependencies are
14 expressed using a set of libraries which bring as much as possible checking and validation
15 to the type level, without sacrificing expressiveness. Together they form a toolkit for defining
16 scalable cloud-based workflows composed of stateless computations, with a static reproducible
17 specification of dependencies, behavior and wiring of all steps. The modeling of taxonomy data is
18 done using Bio4j, where the new paradigm of graph databases allows for both a simple expression
19 of taxonomic assignment tasks and the calculation of taxa abundance values considering the
20 hierarchic structure of the taxonomy tree. MG7 includes a new 16S reference database, *16S-*
21 *DB7*, built with a flexible and sustainable update system, and the possibility of project-driven
22 personalization.

23 † The first and second authors contributed equally to this work.

24 **Keywords:** Metagenomics, 16S, Bacterial diversity profile, Bio4j, Graph databases, Cloud computing, NGS, Genomic big data,
25 Microbiome, Environmental, 16S Database

1 INTRODUCTION

26 During the past decade, metagenomics data analysis is growing exponentially. Some of the reasons behind
27 this are the increasing throughput of massively parallel sequencing technologies (with the derived decrease
28 in sequencing costs), and the wide impact of metagenomics studies (Oulas et al., 2015), especially in
29 human health (diagnostics, treatments, drug response or prevention) (Bikel et al., 2015). We should also
30 mention what could be called the microbiome explosion: all kind of microbiomes (gut, mouth, skin, urinary
31 tract, airway, milk, bladder) are now routinely sequenced in different conditions of health and disease,
32 or after different treatments. The impact of Metagenomics is also being felt in environmental sciences
33 (Ufarté et al., 2015), crop sciences, the agrifood sector (Coughlan et al., 2015) and biotechnology in general
34 (Cowan et al., 2015, Kodzius and Gojobori (2015)). These new possibilities for exploring the diversity of
35 micro-organisms in the most varied environments are opening new research areas, and drastically changing
36 the existing ones.

37 As a consequence, the challenge is thus moving (as in other fields) from data acquisition to data analysis:
38 the amount of data is expected to be overwhelming in a very short time (Stephens et al., 2015).

39 Genome researchers have raised the alarm over big data in the past (Hayden, 2015), but even a more
40 serious challenge might be faced with the metagenomics boom. If we compare metagenomics data with
41 other genomics data used in clinical genotyping we find a differential feature: the key role of time. Thus,
42 for example, in some longitudinal studies, serial sampling from the same patient (Faust et al., 2015) along
43 several weeks (or years) is being used for the follow up of some intestinal pathologies, for studying the
44 evolution of the gut microbiome after antibiotic treatment, or for colon cancer early detection (Zeller et al.,
45 2014, Garrett (2015)). This need of sampling across time adds more complexity to metagenomics data
46 storage and demands adapted algorithms to detect state variations across time as well as idiosyncratic
47 commonalities of the microbiome of each individual (Franzosa et al., 2015). In addition to the intra-
48 individual sampling-time dependence, metagenomic clinical test results vary depending on the specific
49 region of extraction of the clinical specimen. This local variability adds complexity to the analysis since
50 different localizations (different tissues, different anatomical regions, healthy or tumor tissues) are required
51 to have a sufficiently complete landscape of the human microbiome. Moreover, re-analysis of old samples
52 using new tools and better reference databases might be also demanded from time to time.

53 Other disciplines such as astronomy or particle physics have faced the big data challenge before. A key
54 difference is the existence of standards for data processing (Stephens et al., 2015); in metagenomics global
55 standards for converting raw sequence data into processed data are not yet well defined, and there are
56 shortcomings derived from the fact that most bioinformatics methodologies used for metagenomics data
57 analysis were designed for scenarios very different from the current one. These are some of the aspects that
58 have suffered crucial changes and advances with a direct impact in metagenomics data analysis:

- 59 1. **Sequence data:** the reads are larger, the sequencing depth and the number of samples of each project
60 are considerably bigger. The first metagenomics studies were very local projects, while nowadays
61 the most fruitful studies are done at a global level (international, continental, national). This kind of
62 global studies has yielded the discovery of clinical biomarkers for diseases of the importance of cancer,
63 obesity or inflammatory bowel diseases and has allowed exploring the biodiversity of varied earth
64 environments.
- 65 2. **The genomics explosion:** its effect being felt in this case in the reference sequences. The immense
66 amount of sequences available in public repositories demands new strategies for curation, update and

67 storage of metagenomics reference databases: current models will (already) have problems to face the
68 future avalanche of metagenomic sequence data.

69 3. **Cloud computing:** the appearance of new models for massive computation and storage such as the
70 cloud-based platforms, or the widespread adoption of programming methodologies like functional
71 programming, or, more speculatively, dependently typed programming. The new possibilities that
72 these advances offer must have a direct impact in metagenomics data analysis.

73 4. **Open science:** the new social manner to do science, particularly so in genomics, brings its own
74 set of requirements. Metagenomics evolves in a social and global scenario following a science
75 democratization trend in which many small research groups from distant countries share a common
76 big metagenomics project; this global cooperation demands systems allowing for reproducible data
77 analysis, data interoperability, and tools and practices for asynchronous collaboration between different
78 groups.

2 RESULTS

79 2.1 Overview

80 Considering the current new metagenomics scenario and to tackle the challenges posed by metagenomics
81 big data analysis outlined in the Introduction we have designed a new open source methodology for
82 analyzing metagenomics data. It exploits the new possibilities that cloud computing offers to get a system
83 robust, programmatically configurable, modular, distributed, flexible, scalable and traceable in which the
84 biological databases of reference sequences can be easily updated and/or frequently substituted by new
85 ones or by databases specifically designed for focused projects.

86 These are some of the more innovative MG7 features:

- 87 • Static reproducible specification of dependencies and behavior of the different components using
88 *Statika* and *Datasets*
- 89 • Parallelization and distributed analysis based on AWS, with on-demand infrastructure as the basic
90 paradigm
- 91 • Definition of complex workflows using *Loquat*, a composable system for scaling/parallelizing stateless
92 computations especially designed for AWS
- 93 • A new approach to data analysis specification, management and specification based on working with
94 it in exactly the same way as for a software project, together with the extensive use of compile-time
95 structures and checks
- 96 • Modeling of the taxonomy tree using the new paradigm of graph databases (Bio4j). It facilitates the
97 taxonomic assignment tasks and the calculation of the taxa abundance values considering the hierarchic
98 structure of taxonomy tree (cumulative values)
- 99 • Exhaustive per-read taxonomic assignment using two complementary assignment algorithms Lowest
100 Common Ancestor and Best BLAST Hit
- 101 • Using a new 16S database of reference sequences (16S-DB7) with a flexible and sustainable system of
102 updating and project-driven customization

103 2.2 Libraries and resources

104 In this section we describe the resources and libraries developed by the authors on top of which MG7
105 is built. All MG7 code is written in Scala, a hybrid object-functional programming language. Scala was

106 chosen based on the possibility of using certain advanced programming styles, and Java interoperability,
107 which let us build on the vast number of existing Java libraries; we take advantage of this when using
108 Bio4j as an API for the NCBI taxonomy. It has support for type-level programming, type-dependent types
109 (through type members) and singleton types, which permits a restricted form of dependent types where
110 types can depend essentially on values determined at compile time (through their corresponding singleton
111 types). Conversely, through implicits one can retrieve the value corresponding to a singleton type.

112 2.2.1 *Statika*: machine configuration and behavior

113 *Statika* is a Scala library developed by the first and last authors which serves as a way of defining
114 and composing machine behaviors statically. The main component are **bundles**. Each bundle declares a
115 sequence of computations (its behavior) which will be executed in an **environment**. A bundle can *depend*
116 on other bundles, and when being executed by an environment, its DAG (Directed Acyclic Graph) of
117 dependencies is linearized and run in sequence. In our use, bundles correspond to what an EC2 instance
118 should do and an environment to an AMI (Amazon Machine Image) which prepares the basic configuration,
119 downloads the Scala code and runs it.

120 2.2.2 *Datasets*: a mini-language for data

121 *Datasets* is a Scala library developed by the first and last authors with the goal of being a Scala-embedded
122 mini-language for datasets and their locations. **Data** is represented as type-indexed fields: keys are modeled
123 as singleton types, and values correspond to what could be called a denotation of the key: a value of
124 type `Location` tagged with the key type. Then a **Dataset** is essentially a collection of data, which
125 are guaranteed statically to be different through type-level predicates, making use of the value–type
126 correspondence which can be established through singleton types and implicits. A dataset location is then
127 just a list of locations formed by locations of each dataset key. All this is based on what could be described
128 as an embedding in Scala of an extensible record system with concatenation on disjoint labels, in the
129 spirit of (Harper and Pierce, 1990, Harper and Pierce (1991)). For that *Datasets* uses *ohnosequences/cosas*
130 library.

131 Data keys can further have a reference to a **data type**, which, as the name hints at, can help in providing
132 information about the type of data we are working with. For example, when declaring Illumina reads as a
133 data, a data type containing information about the read length, insert size or end type (single or paired) is
134 used.

135 A **location** can be, for example, an S3 object or a local file; by leaving the location type used to denote
136 particular data free we can work with different “physical” representations, while keeping track of to which
137 logical data they are a representation of. Thus, a process can generate locally a `.fastq` file representing
138 the merged reads, while another can put it in S3 with the fact that they all correspond to the “same” merged
139 reads is always present, as the data that those “physical” representations denote.

140 2.2.3 *Loquat*: Parallel data processing with AWS

141 *Loquat* is a library developed by the first, second and last authors designed for the execution of
142 embarrassingly parallel tasks using S3, SQS and EC2 Amazon services.

143 A *loquat* executes a process with explicit input and output datasets (declared using the *Datasets* library
144 described above). Workers (EC2 instances) read from an SQS queue the S3 locations for both input and
145 output data; then they download the input to local files, and pass these file locations to the process to be
146 executed. The output is then put in the corresponding S3 locations.

147 A manager instance is used to monitor workers, provide initial data to be put in the SQS queue and
148 optionally release resources depending on a set of configurable conditions.

149 Both worker and manager instances are *Statika* bundles. The worker can declare any dependencies needed
150 to perform its task: other tools, libraries, or data.

151 All configuration such as the number of workers or the instance types is declared statically, the
152 specification of a loquat being ultimately a Scala object. Deploy and resource management methods
153 make easy to use an existing loquat either as a library or from (for example) a Scala REPL.

154 The input and output (and their locations) being defined statically has several critical advantages. First,
155 composing different loquats is easy and safe; just use the output types and locations of the first one as input
156 for the second one. Second, data and their types help in not mixing different resources when implementing
157 a process, while serving as a safe and convenient mechanism for writing generic processing tasks. For
158 example, merging paired-end Illumina reads generically is easy as the data type includes the relevant
159 information (insert size, read length, etc) to pass to a tool such as FLASH.

160 2.2.4 Type-safe eDSLs for BLAST and FLASH

161 We developed our own Scala-based type-safe eDSLs (embedded Domain Specific Languages) for FLASH
162 (Magoč and Salzberg, 2011) and BLAST (Camacho et al., 2009) expressions and their execution.

163 In the case of BLAST we use a model where we can guarantee for each BLAST command expression at
164 compile time that

- 165 • all required arguments are provided
- 166 • only valid options are provided
- 167 • correct types for each option value
- 168 • valid output record specification

169 Generic type-safe parsers returning a heterogeneous record of BLAST output fields are also available,
170 together with output data defined using *Datasets* which have a reference to the exact BLAST command
171 options which yielded that output. This lets us provide generic parsers for BLAST output which are
172 guaranteed to be correct.

173 In the same spirit as for BLAST, we implemented a type-safe eDSL for FLASH expressions and their
174 execution, supporting features equivalent to those outlined for the BLAST eDSL.

175 2.2.5 Bio4j and Graph Databases

176 (Bio4j Pareja-Tobes et al., 2015) is a data platform integrating data from different resources such as
177 UniProt or GO in a graph data paradigm. In the assignment phase we use a subgraph containing the NCBI
178 Taxonomy, wrapping in Scala its Java API in a tree algebraic data type.

179 2.2.6 16S-DB7 Reference Database Construction

180 Our 16S-DB7 Reference Database is a curated subset of sequences from the NCBI nucleotide database
181 **nt**. The sequences included were selected by similarity with the bacterial and archaeal reference sequences
182 downloaded from the **RDP database** (Cole et al., 2013). RDP unaligned sequences were used to capture
183 new 16S RNA sequences from **nt** using BLAST similarity search strategies and then, performing
184 additional curation steps to remove sequences with poor taxonomic assignments to taxonomic nodes
185 close to the root of the taxonomy tree. All the nucleotide sequences included in **nt** database has a

186 taxonomic assignment provided by the **Genbank** sequence submitter. NCBI provides a table (available
187 at <ftp://ftp.ncbi.nlm.nih.gov/pub/taxonomy/>) to do the mapping of any Genbank Identifier (GI) to its
188 Taxonomy Identifier (TaxID). Thus, we are based on a crowdsourced submitter-maintained taxonomic
189 annotation system for reference sequences. It supposes a sustainable system able to face the expected
190 number of reference sequences that will populate the public global nucleotide databases in the near future.
191 Another advantageous point is that we are based on NCBI taxonomy, the *de-facto* standard taxonomic
192 classification for biomolecular data (Cochrane and Galperin, 2010). NCBI taxonomy is, undoubtedly, the
193 most used taxonomy all over the world and the most similar to the official taxonomies of each specific field.
194 This is a crucial point because all the type-culture and tissue databanks follow this official taxonomical
195 classification and, in addition, all the knowledge accumulated during last decades is referred to this
196 taxonomy. In addition NCBI provides a direct connection between taxonomical formal names and the
197 physical specimens that serve as exemplars for the species (Federhen, 2014).

198 Certainly, if metagenomics results are easily integrated with the theoretical and experimental knowledge
199 of each specific area, the impact of metagenomics will be higher than if it progresses as a disconnected
200 research branch. Considering that metagenomics data interoperability, which is especially critical in clinical
201 environments, requires a stable taxonomy to be used as reference, we decided to rely on the most widely
202 used taxonomy: the NCBI taxonomy. In addition, the biggest global sequence database GenBank follows
203 this taxonomy to register the origin of all their submitted sequences. Our 16S database building strategy
204 allows the substitution of the 16S database by any other subset of **nt**, even by the complete **nt** database if it
205 would be needed, for example, for analyzing shotgun metagenomics data. This possibility of changing
206 the reference database provides flexibility to the system enabling it for easy updating and project-driven
207 personalization.

208 **2.3 Workflow Description**

209 The MG7 analysis workflow is summarized in Figure 1. The input files for MG7 are the FASTQ files
210 resulting from a paired-end NGS sequencing experiment.

211 **2.3.1 Joining reads of each pair using FLASH**

212 In the first step the paired-end reads, designed with an insert size that yields pairs of reads with an
213 overlapping region between them, are assembled using FLASH (Magoč and Salzberg, 2011). FLASH is
214 designed to merge pairs of reads when the original DNA fragments are shorter than twice the length of
215 reads. Thus, the sequence obtained after joining the 2 reads of each pair is larger and has better quality
216 since the sequence at the ends of the reads is refined merging both ends in the assembly. To have a larger
217 and improved sequence is crucial to do more precise the inference of the bacterial origin based on similarity
218 with reference sequences.

219 **2.3.2 Parallelized BLASTN of each read against the 16S-DB7**

220 The second step is to search for similar 16S sequences in our 16S-DB7 database. The taxonomic
221 assignment for each read is based on BLASTN of each read against the 16S database. Assignment based on
222 direct similarity of each read one by one compared against a sufficiently wide database is considered in
223 different reviews of metagenomics analysis methodologies (Segata et al., 2013, Morgan and Huttenhower
224 (2012)) as a very exhaustive method for assignment. Some methods of assignment compare the sequences
225 only against the 16S genes from available complete bacterial genomes or avoid computational cost clustering
226 or binning the sequences first, and then doing the assignments only for the representative sequence of each
227 cluster. MG7 carries out an exhaustive comparison of all the reads under analysis and it does not applies

228 any binning strategy. Every read is specifically compared with all the sequences of the 16S database. We
229 select the best BLAST hits (10 hits by default) obtained for each read to do the taxonomic assignment.

230 2.3.3 Taxonomic Assignment Algorithms

231 All the reads are assigned under two different algorithms of assignment: i. Lowest Common Ancestor
232 based taxonomic assignment (LCA) and ii. Best BLAST Hit based taxonomic assignment (BBH). Figure 2
233 displays schematically the LCA algorithm applied *sensu stricto* (left panel) and the called ‘in line’ exception
234 (right panel) designed in order to gain specificity in the assignments in the cases in which the topology of
235 the taxonomical nodes corresponding to the BLAST hits support sufficiently the assignment to the most
236 specific taxon.

237 2.3.3.1 Lowest Common Ancestor based Taxonomic Assignment

238 For each read, first, we select the BEST BLAST HITs (by default 10 Hits) over a threshold of similarity.
239 To evaluate similarity for this first filtering of hits we use the Expect value (by default $evaluate \leq e^{-15}$) that
240 describes the number of hits one can “expect” to see by chance when searching a database of a particular
241 size. In a second filtering step we filtering those hits that are not sufficiently good comparing them with the
242 best one. We select the best HSP (High Similarity Pair) per reference sequence and then choose the best
243 HSP (that with lowest E-value) between all the selected ones. The bitscore of this best HSP (called S) is
244 used as reference to filter the rest of HSPs. All the HSPs with bitscore below the product pS are filtered. p
245 is a coefficient fixed by the user to define the bitscore required, e.g. if $p = 0.9$ and $S = 700$ the required
246 bitscore threshold would be 630. Once we have the definitive HSPs selected, we obtain their corresponding
247 taxonomic nodes using the taxonomic assignments that NCBI provides for all the nt database sequences.
248 Now we have to analyze the topological distribution of these nodes in the taxonomy tree: i. If all the nodes
249 forms a line in the taxonomy tree (are located in a not branched lineage to the tree root) we should choose
250 the most specific taxID as the final assignment for that read. We call to this kind of assignment the ‘in line’
251 exception (see Figure 2 right panel). ii. If not, we should search for the *sensu stricto* Lowest Common
252 Ancestor (LCA) of all the selected taxonomic nodes (See Figure 2 left panel). In this approach we decided
253 to use the bitscore for evaluating the similarity because it is a value that increases when similarity is higher
254 and depends a lot on the length of the HSP. Some reads could not find sequences with enough similarity in
255 the database and then they would be classified as reads with no hits. Advanced metagenomics analysis
256 approaches (Huson and Weber, 2012) have adopted LCA assignment algorithms because it provides fine
257 and trusted taxonomical assignment.

258 2.3.3.2 Best BLAST hit taxonomic assignment

259 We decided to maintain the simpler method of Best BLAST Hit (BBH) for taxonomic assignment because,
260 in some cases, it can provide information about the sequences that adds information to that obtained using
261 LCA algorithm. Using LCA algorithm, when some reference sequences with BLAST alignments over
262 the required thresholds map to a not sufficiently specific taxID, the read can be assigned to an unspecific
263 taxon near to the root of the taxonomy tree. If the BBH reference sequence maps to more specific taxa, this
264 method, in that case, gives us useful information.

265 2.3.4 Output for LCA and BBH assignments

266 MG7 provides independent results for the 2 different approaches, LCA and BBH. The output files include,
267 for each taxonomy node (with some read assigned), abundance values for direct assignment and cumulative
268 assignment. The abundances are provided in counts (absolute values) and in percentage normalized to the
269 number of reads of each sample. Direct assignments are calculated counting reads specifically assigned to a

270 taxonomic node, not including the reads assigned to the descendant nodes in the taxonomy tree. Cumulative
271 assignments are calculated including the direct assignments and also the assignments of the descendant
272 nodes. For each sample MG7 provides 8 kinds of abundance values: LCA direct counts, LCA cumu. counts,
273 LCA direct %, LCA cumu. %, BBH direct counts, BBH cumu. counts, BBH direct % and BBH cumu. %.

274 2.4 Data analysis as a software project

275 The MG7 16S data analysis workflow is indeed a set of tasks, all of them based in *Loquat*. For each task,
276 a set of inputs and outputs as well as configuration parameters must be statically defined. The user is also
277 free to leave the reasonable defaults for configuration, needing only to define the input and output of the
278 whole workflow. The definition of this configuration is Scala code and the way of starting an MG7 analysis
279 is compiling the project code and launching it from the Scala interactive console.

280 Code compilation prior to launching any analysis assures that no AWS resources are launched if the
281 analysis is not well-defined, avoiding expenses not leading to any analysis. Besides compile-time checks,
282 runtime checks are made before launch to ensure existence of input data and availability of resources.

283 An MG7 analysis is then a Scala project where the user only needs to set certain variables at the code
284 level (input, output and parameters), compile the code and run it. To facilitate the process of setting up the
285 Scala project, a template with sensible defaults is provided.

286 In order to be able to exploit AWS infrastructure for the MG7 analysis, the user needs to set up an AWS
287 account with certain IAM (Identity and Access Management) permission policies that will grant access to
288 the resources used in the workflow.

289 2.5 Availability

290 MG7 is open source, available at <https://github.com/ohnosequences/mg7> under an AGPLv3 license.

3 DISCUSSION

291 We could summarize the most innovative ideas and developments in MG7:

- 292 1. Treating data analysis as a software project. This makes for radical improvements in *reproducibility*,
293 *reuse*, *versioning*, *safety*, *automation* and *expressiveness*
- 294 2. Checking at compile-time: input and output data, their locations and type are expressible and checked
295 at compile-time using *Datasets*
- 296 3. Management of dependencies and machine configurations using *Statika*
- 297 4. Automation of AWS cloud resources and processes, including distribution and parallelization through
298 the use of *Loquat*
- 299 5. Taxonomic data and related operations are treated natively as what they are: graphs, through the use of
300 *Bio4j*
- 301 6. MG7 provides a sustainable model for taxonomic assignment, appropriate to face the challenging
302 amount of data that high throughput sequencing technologies generate

303 We will expand on each item in the following sections.

304 3.1 A new approach to data analysis: data analysis as a software project and checking 305 at compile-time

306 MG7 proposes to define and work with a particular data analysis task as a software project, using Scala.
307 The idea is that *everything*: data description, their location, configuration parameters and the infrastructure
308 used should be expressed as Scala code, and treated in the same way as any (well-managed) software
309 project. This includes, among other things, using version control systems (`git` in our case), writing tests,
310 making stable releases following semantic versioning or publishing artifacts to a repository.

311 What we see as key advantages of this approach (when coupled with compile-time specification and
312 checking), are

- 313 • **Reproducibility** the same analysis can be run again with exactly the same configuration in a trivial
314 way.
- 315 • **Versioning** as in any software project, there can be different versions, stable releases, etc.
- 316 • **Reuse** we can build standard configurations on top of this and reuse them for subsequent data analysis.
317 A particular data analysis *task* can be used as a *library* in further analysis.
- 318 • **Decoupling** We can start working on the analysis specification, without any need for available data in
319 a much easier way.
- 320 • **Documentation** We can take advantage of all the effort put into software documentation tools and
321 practices, such as in our case Scaladoc or literate programming. As documentation, analysis processes
322 and data specification live together in the files, it is much easier to keep coherence between them.
- 323 • **Expresiveness and safety** For example in our case we can choose only from valid illumina read types,
324 and then build a default FLASH command based on that. The output locations, being declared statically,
325 are also available for use in further analysis.

326 3.2 Input and output data declaration

327 An important aspect of the MG7 workflow is the way it deals with data resources. All the data that is
328 going to be used in the analysis or produced as an output is described as Scala code using rich types from
329 the *Datasets* language. This allows the user to specify information about types of data, information that can
330 then be utilized by tools analyzing this data. For example, we can specify that, for the first part of the MG7
331 workflow, running FLASH in parallel requires illumina paired end reads and produces joined reads.

332 On one hand, specification of the input data allows us to restrict its type and force users to be conscious
333 about what they pass as an input. On the other hand, specification of the output data helps to build a
334 workflow as a *composition* of several parts: we can ensure on the Scala code type level that the output
335 of one component fits as an input for the next component. This is crucial as, obviously, the way a data
336 analysis task works depends a lot on the particular structure of the data. For instance, in the MG7 workflow,
337 using BLAST eDSL, we can precisely describe which format will have the output of the BLAST step,
338 which information it will include, and then in the next step we can reuse this description to parse BLAST
339 output and retrieve the part of the information needed for the taxonomy assignment analysis. Having the
340 data structure described statically as Scala code allows us to be sure that we will not have parsing problems
341 or other issues with incompatible data passed between workflow components.

342 All this does not compromise flexibility in how the user works with data in MG7: having static data
343 declarations as a part of the configuration allows the user to reuse analysis components, or modify them
344 according to particular needs. Besides that, an important advantage of the type-level control is the added

345 protection from the execution (and deployment) of a wrongly configured analysis task, which may lead to
346 significant costs in both time and money.

347 **3.3 Tools, data, dependencies and automated deployment**

348 Bioinformatics software often has a complicated installation process and requires various dependencies
349 with unclear versions. This makes the deployment of the bioinformatics tools an involved task and resolving
350 it manually is not a solution in the context of cloud computations. To face this problem, one needs an
351 automated system of managing tools and resources, which will allow an expressive way for describing
352 dependencies between parts of a pipeline and provide a reproducible procedure of its deployment. We have
353 developed *Statika* for this purpose and successfully used it in MG7.

354 Every external tool involved in the workflow is represented as a *Statika* bundle, which is essentially a
355 Scala project describing the installation process of this tool and declaring dependencies on other bundles
356 which will be installed prior to the considered tool itself. Describing relationships between bundles on
357 the code level allows us to track the directed acyclic graph of their dependencies and linearize them to
358 automatically install them sequentially in the right order. Meanwhile, describing the installation process
359 on the code level allows the user to utilize the wide range of available Scala and Java APIs and tools,
360 making installation a well-defined sequence of steps rather than an unreliable script, dependent on a certain
361 environment. *Statika* offers an easy path towards making deployment an automated, reproducible process.

362 Besides bioinformatics tools like BLAST and FLASH, *Statika* bundles are used for wrapping data
363 dependencies and all inner components of the system that require cloud deployment. In particular, all
364 components of *Loquat* are bundles; the user can then define which components are needed for the parallel
365 processing on each computation unit in an expressive way, declaring them as bundle dependencies of the
366 loquat “worker” bundle. This modularization is also important for the matter of making components of
367 the system reusable for different projects and liberating the user from most of the tasks related to their
368 deployment.

369 **3.4 Parallel computations in the cloud**

370 The MG7 workflow consists of certain steps, each of which performs some work in parallel, using
371 the cloud infrastructure managed by *Loquat*. It is important to notice the horizontal scalability of this
372 approach. Irrespectively of how much data needs to be processed, MG7 will handle it, by splitting data
373 into chunks and performing the analysis on multiple computation units. The Amazon Elastic Compute
374 Cloud (EC2) service provides a transparent way of managing computation infrastructure, called autoscaling
375 groups. The User can set MG7 configuration parameters, adjusting for each task the amount and hardware
376 characteristics of the EC2 instances they want to use for it. But it is important to note that, as each workflow
377 step is not very resource demanding, it is not needed to hire EC2 instances with some advanced hardware.
378 Instead, an average type will work and you can reduce execution time by simply scaling out the number of
379 instances.

380 **3.5 Taxonomy and Bio4j**

381 The hierarchic structure of the taxonomy of the living organisms is a tree, and, hence, is also a graph
382 in which each node, with the exception of the root node, has a unique parent node. It led us to model the
383 taxonomy tree as a graph using the graph database paradigm. Previously we developed Bio4j [Pareja-
384 Tobes-2015], a platform for the integration of semantically rich biological data using typed graph models.
385 It integrates most publicly available data linked with sequences into a set of interdependent graphs to

386 be used for bioinformatics analysis and especially for biological data. MG7 works based on the Bio4j
387 taxonomy module. It opens the possibility to connect the taxonomic profiling data obtained with MG7 to
388 all the biological knowledge associated to each taxon. Using the information available in Bio4j for all the
389 proteins assigned to each taxon we are connected to all the functional data available in Uniprot related with
390 it.

391 **3.6 Future developments**

392 3.6.1 Shotgun metagenomics

393 It is certainly possible to adapt MG7 to work with shotgun metagenomics data. Simply changing the
394 reference database to include whole genome sequence data could yield interesting results. This could also
395 be refined by restricting reference sequences according to all sort of criteria, like biological function or
396 taxonomy. Bio4j would be an invaluable tool here, thanks to its ability to express complex predicates on
397 sequences using all the information linked with them (GO annotations, UniProt data, NCBI taxonomy, etc).

398 3.6.2 Comparing groups of samples

399 The comparison of the taxonomic profiles between different groups of samples is a need for many
400 metagenomics studies. Tasks related with this group-based analysis, such as the extraction of the minimal
401 tree with all the taxa with some direct or accumulated assignment, will be part of a new MG7 module,
402 already in development.

403 3.6.3 Interactive visualizations based on Biographika

404 New visualization tools for metagenomics results are undoubtedly needed. Interactivity is a especially
405 interesting feature for metagenomics data visualization, since the expert needs to explore the results in a
406 knowledge-driven way. The majority of the available metagenomics data visualizations are static. We are
407 working in the *Biographika* project (Tobes et al., 2015), to provide interactive rich visualizations on the
408 web for Bio4j data. The development of visualizations specific for MG7 is one of Biographika current
409 goals. Biographika is based on D3.js, the de-facto standard JavaScript data visualization library, and is
410 open source.

4 MATERIALS AND METHODS

411 **4.1 Amazon Web Services**

412 MG7 uses the following Amazon Web Services:

- 413 • EC2 (Elastic Compute Cloud) autoscaling groups for launching and managing computation units
- 414 • S3 (Simple Storage Service) for storing input and output data
- 415 • SQS (Simple Queue Service) for communication between different components of the system
- 416 • SNS (Simple Notification Service) for e-mail notifications

417 These services are used through a Scala wrapper of the official AWS Java SDK v1.9.25:
418 ohnosequences/aws-scala-tools v0.13.2.

419 **4.2 Scala**

420 MG7 itself and all the libraries used are written in Scala v2.11.

421 4.3 Statika

422 MG7 uses ohnosequences/statika v2.0.0 for specifying the configuration and behavior of EC2 instances.

423 4.4 Datasets

424 MG7 uses ohnosequences/datasets v0.2.0 for specifying input and output data, their type and their
425 location.

426 4.5 Loquat

427 MG7 uses ohnosequences/loquat v2.0.0 for the specification of data processing tasks and their execution
428 using AWS resources.

429 4.6 BLAST eDSL

430 MG7 uses ohnosequences/blast v0.2.0. The BLAST version used is v2.2.31+.

431 4.7 FLASH eDSL

432 MG7 uses ohnosequences/flash v0.1.0. The FLASH version used is v1.2.11.

433 4.8 Bio4j

434 MG7 uses bio4j/bio4j v0.12.0-RC3 and bio4j/bio4j-titan v0.4.0-RC2 as an API for the NCBI taxonomy.

5 DISCLOSURE/CONFLICT-OF-INTEREST STATEMENT

435 All authors work at the *Oh no sequences!* research group, part of Era7 Bioinformatics. Era7 offers
436 metagenomics data analysis services based on MG7. MG7 is open source, available under the OSI-approved
437 AGPLv3 license.

6 AUTHOR CONTRIBUTIONS

- 438 • **AA** developed *MG7*, *Statika*, *Datasets*, and *aws-scala-tools*; wrote the paper;
- 439 • **EK** developed *nispero* (a prototype for *Loquat* (Kovach et al., 2014)) and *aws-scala-tools*.
- 440 • **MM** *MG7* workflow design; curation and design of the *16S-DB7* reference database; wrote the paper.
- 441 • **PPT** curation, design, data extraction code of the *16S-DB7* reference database.
- 442 • **EP** *MG7* workflow design; wrote the paper.
- 443 • **RT** *MG7* workflow design, assignment strategy; curation and design of the *16S-DB7* reference
444 database; wrote the paper.
- 445 • **EPT** developed *MG7*, *Statika*, *Datasets*, *FLASH/BLAST eDSLs*; data analysis approach and design;
446 wrote the paper.

447 All authors have read and approved the final manuscript.

7 ACKNOWLEDGEMENTS

448 *Funding*: The two first authors are funded by INTERCROSSING (Grant 289974).

8 FIGURE 1: MG7 ANALYSIS WORKFLOW.

449 The paired reads in fastq format are merged resulting in only one sequence per read pair. The next step
450 is a parallelized BLASTN of every merged sequence against the 16S reference database 16S-DB7. Then,
451 the mapping of the detected similar sequences in the database to the taxonomy node to which they belong
452 is carried out. This is done using Bio4j that includes a module with all the NCBI taxonomy in a graph
453 connected with the Gene Ontology, Uniprot, and RefSeq graphs. Then the taxonomic assignment is
454 done for each sequence following two different approaches: LCA and BBH, and finally the abundances
455 corresponding to direct and cumulative assignments for each node in percentage and absolute counts are
456 provided for each assignment mode.

9 FIGURE 2: LOWEST COMMON ANCESTOR ALGORITHM FOR TAXONOMIC ASSIGNMENT.

457 The Left panel displays an example of the application of LCA algorithm in a *sensu stricto* mode. A, B,
458 C and D represent taxonomy tree nodes with assigned reads. Right panel displays the *in line* mode of
459 assignment which is an exception for the *sensu stricto* mode of application of LCA algorithm. The *in line*
460 mode is used when all the nodes are located in a line without bifurcations. In that case the taxon assigned
461 is the most specific (the most distant from the root).

REFERENCES

- 462 Bikel, S., Valdez-Lara, A., Cornejo-Granados, F., Rico, K., Canizales-Quinteros, S., Soberón, X., et al.
463 (2015). Combining metagenomics, metatranscriptomics and viromics to explore novel microbial
464 interactions: towards a systems-level understanding of human microbiome. *Computational and structural*
465 *biotechnology journal* 13, 390–401
- 466 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., et al. (2009). Blast+:
467 architecture and applications. *BMC bioinformatics* 10, 421
- 468 Cochrane, G. R. and Galperin, M. Y. (2010). The 2010 nucleic acids research database issue and online
469 database collection: a community of data resources. *Nucleic acids research* 38, D1–D4
- 470 Cole, J. R., Wang, Q., Fish, J. A., Chai, B., McGarrell, D. M., Sun, Y., et al. (2013). Ribosomal database
471 project: data and tools for high throughput rRNA analysis. *Nucleic acids research* , gkt1244
- 472 Coughlan, L. M., Cotter, P. D., Hill, C., and Alvarez-Ordóñez, A. (2015). Biotechnological applications of
473 functional metagenomics in the food and pharmaceutical industries. *Frontiers in microbiology* 6
- 474 Cowan, D. A., Ramond, J.-B., Makhalanyane, T. P., and De Maayer, P. (2015). Metagenomics of extreme
475 environments. *Current opinion in microbiology* 25, 97–102
- 476 Faust, K., Lahti, L., Gonze, D., de Vos, W. M., and Raes, J. (2015). Metagenomics meets time series
477 analysis: unraveling microbial community dynamics. *Current opinion in microbiology* 25, 56–66
- 478 Federhen, S. (2014). Type material in the ncbi taxonomy database. *Nucleic acids research* , gku1127
- 479 Franzosa, E. A., Huang, K., Meadow, J. F., Gevers, D., Lemon, K. P., Bohannon, B. J., et al. (2015).
480 Identifying personal microbiomes using metagenomic codes. *Proceedings of the National Academy of*
481 *Sciences* , 201423854
- 482 Garrett, W. S. (2015). Cancer and the microbiota. *Science* 348, 80–86
- 483 Harper, R. and Pierce, B. (1991). A record calculus based on symmetric concatenation. In *Proceedings of*
484 *the 18th ACM SIGPLAN-SIGACT symposium on Principles of programming languages* (ACM), 131–142
- 485 Harper, R. W. and Pierce, B. C. (1990). Extensible records without subsumption

- 486 Hayden, E. C. (2015). Genome researchers raise alarm over big data. *Nature*
- 487 Huson, D. H. and Weber, N. (2012). Microbial community analysis using megan. *Methods in enzymology*
- 488 531, 465–485
- 489 Kodzius, R. and Gojobori, T. (2015). Marine metagenomics as a source for bioprospecting. *Marine*
- 490 *genomics*
- 491 Kovach, E., Alekhin, A., Manrique, M., Pareja-Tobes, P., Pareja, E., Tobes, R., et al. (2014). Nispero:
- 492 a cloud-computing based scala tool specially suited for bioinformatics data processing. In *IWBBIO*.
- 493 1414–1415
- 494 Magoč, T. and Salzberg, S. L. (2011). Flash: fast length adjustment of short reads to improve genome
- 495 assemblies. *Bioinformatics* 27, 2957–2963
- 496 Morgan, X. C. and Huttenhower, C. (2012). Chapter 12: human microbiome analysis. *PLoS Comput Biol*
- 497 8, e1002808
- 498 Oulas, A., Pavlodi, C., Polymenakou, P., Pavlopoulos, G. A., Papanikolaou, N., Kotoulas, G., et al.
- 499 (2015). Metagenomics: Tools and insights for analyzing next-generation sequencing data derived from
- 500 biodiversity studies. *Bioinformatics and biology insights* 9, 75
- 501 Pareja-Tobes, P., Tobes, R., Manrique, M., Pareja, E., and Pareja-Tobes, E. (2015). Bio4j: a high-
- 502 performance cloud-enabled graph-based data platform. *bioRxiv* , 016758
- 503 Segata, N., Boernigen, D., Tickle, T. L., Morgan, X. C., Garrett, W. S., and Huttenhower, C. (2013).
- 504 Computational meta’omics for microbial community studies. *Molecular systems biology* 9, 666
- 505 Stephens, Z. D., Lee, S. Y., Faghri, F., Campbell, R. H., Zhai, C., Efron, M. J., et al. (2015). Big data:
- 506 Astronomical or genetical? *PLoS Biol* 13, e1002195
- 507 Tobes, P. P., Tobes, E. P., Manrique, M., Pareja, E., and Tobes, R. (2015). Biographika: rich interactive
- 508 data visualizations on the web for the research community. *bioRxiv* , 021063
- 509 Ufarté, L., Potocki-Véronèse, G., and Laville, E. (2015). Discovery of new protein families and functions:
- 510 new challenges in functional metagenomics for biotechnologies and microbial ecology. *Name: Frontiers*
- 511 *in Microbiology* 6, 563
- 512 Zeller, G., Tap, J., Voigt, A. Y., Sunagawa, S., Kultima, J. R., Costea, P. I., et al. (2014). Potential of fecal
- 513 microbiota for early-stage detection of colorectal cancer. *Molecular systems biology* 10, 766