

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

Alexey Alekhin<sup>1</sup> Evdokim Kovach<sup>1</sup> Marina Manrique<sup>1</sup> Pablo Pareja<sup>1</sup> Eduardo Pareja<sup>1</sup> Raquel Tobes<sup>1</sup> and Eduardo Pareja-Tobes<sup>1,\*</sup>

<sup>1</sup>*Oh no sequences! Research Group, Era7 Bioinformatics, Granada, Spain*

Correspondence\*:

Corresponding Author

Oh no sequences! Research Group, Era7 Bioinformatics, Plaza Campo Verde 3, Granada, 18001, Spain, [eparejatobes@ohnosequences.com](mailto:eparejatobes@ohnosequences.com)

## 2 ABSTRACT

3 No abstract yet. Will be here.

4 **Keywords:** Metagenomics, 16S, Bacterial diversity profile, Bio4j, Graph databases, Cloud computing, NGS, Genomic big data

## 1 1. INTRODUCTION

5 Metagenomics data analysis is growing at exponential rate during the last years. The increasing throughput  
6 of massively parallel sequencing technologies, the derived decreasing cost, and the high impact of  
7 metagenomics studies, especially in human health (diagnostics, treatments, drug response, prevention), are  
8 crucial reasons responsible for this growth of Metagenomics. There is a growing interest in sequencing  
9 all kind of microbiomes (gut, mouth, skin, urinary tract, airway, milk, bladder), in different conditions of  
10 health and disease, or after different treatments. Metagenomics is also impacting environmental sciences,  
11 crop sciences, agrifood sector and biotechnology in general. This new possibilities for exploring the  
12 diversity of micro-organisms in the most diverse environments is opening many new research areas but,  
13 due to this wide interest, it is expected that the amount of data will be overwhelming in the short time (?).

14 Genome researchers have raised the alarm over big data in the past nature news add ref but even a more  
15 serious challenge might be faced with the metagenomics boom/ upswing. If we compare metagenomics  
16 data with other genomics data used in clinical genotyping we find a differential feature: the key role of  
17 time. Thus, for example, in some longitudinal studies, serial sampling of the same patient along several  
18 weeks (or years) is being used for the follow up of some intestinal pathologies, for studying the evolution  
19 of gut microbiome after antibiotic treatment, or for colon cancer early detection (?). This need of sampling  
20 across time adds more complexity to metagenomics data storage and demands adapted algorithms to  
21 detect state variations across time as well as idiosyncratic commonalities of the microbiome of each  
22 individual (?). In addition to the intra-individual sampling-time dependence, metagenomic clinical test  
23 results vary depending on the specific region of extraction of the clinical specimen. This local variability  
24 adds complexity to the analysis since different localizations (different tissues, different anatomical regions,  
25 healthy or tumour tissues) are required to have a sufficiently complete landscape of the human microbiome.

26 Moreover, reanalysis of old samples using new tools and better reference databases might be also demanded  
27 from time to time.

28 During the last years other sciences as astronomy or particle physics are facing the big data challenge but,  
29 at least, these science have standards for data processing (?). Global standards for converting raw sequence  
30 data into processed data are not yet well defined in metagenomics and there are shortcomings derived  
31 from the fact that many bioinformatics methodologies currently used for metagenomics data analysis were  
32 designed for a scenario very different that the current one. These are some of the aspects that have suffered  
33 crucial changes and advances with a direct impact in metagenomics data analysis. i. The first aspect is  
34 related to the sequences to be analyzed: the reads are larger, the sequencing depth and the number of  
35 samples of each project are considerably bigger. The first metagenomics studies were very local projects,  
36 while nowadays the most fruitful studies are done at a global level (international, continental, national).  
37 This kind of global studies has yielded the discovery of clinical biomarkers for diseases of the importance  
38 of cancer, obesity or inflammatory bowel diseases and has allowed exploring the biodiversity in many earth  
39 environments ii. The second aspect derives from the impressive genomics explosion, its effect being felt in  
40 this case in the reference sequences. The immense amount of sequences available in public repositories  
41 demands new approaches in curation, update and storage for metagenomics reference databases: current  
42 models will or already have problems to face the future avalanche of metagenomic sequences. iii. The third  
43 aspect to consider for metagenomics data analysis is related to the appearance of new models for massive  
44 computation and storage and to the new programming methodologies (Scala, ...) and new cloud models  
45 and resources. The immense new possibilities that these advances offer must have a direct impact in the  
46 metagenomics data analysis. iv. And finally the new social manner to do science, and especially genomic  
47 science is the fourth aspect to consider. Metagenomics evolves in a social and global scenario following a  
48 science democratization trend in which many small research groups from distant countries share a common  
49 big metagenomics project. This global cooperation demands systems allowing following exactly the same  
50 pipelines using equivalent cloud resources to modularly execute the analysis in an asynchronous way of  
51 working between different groups. This definitively new scenario demands new methods and tools to handle  
52 the current and future volume of metagenomic data with the sufficient speed of analysis. Considering  
53 all these aspects we have designed a new open source methodology for analyzing metagenomics data  
54 that exploits the new possibilities that cloud computing offers to get a system robust, programmatically  
55 configurable, modular, distributed, flexible, scalable and traceable in which the biological databases of  
56 reference sequences can be easily updated and/or frequently substituted by new ones or by databases  
57 specifically designed for focused projects.

## 2 2. MATERIALS AND METHODS

### 58 2.1 2.x Amazon Web Services

### 59 2.2 2.x Scala

60 Scala is a hybrid object-functional programming language which runs on Java Virtual Machine. It  
61 has support for type-level programming, type-dependent types (through type members) and singleton  
62 types, which permits a restricted form of dependent types where types can depend essentially on values  
63 determined at compile time (through their corresponding singleton types). Conversely, through implicits  
64 one can retrieve the value corresponding to a singleton type.

65 The other key feature for us is Java interoperability, which let us build on the vast number of existing  
66 Java libraries; we take advantage of this when using Bio4j as an API for the NCBI taxonomy.

67 MG7 itself and all the libraries used are written in Scala 2.11.

## 68 2.3 2.x Statika

69 Statika is a Scala library developed by **so and so** which serves as a way of defining and composing  
70 machine behaviors statically. The main component are **bundles**. Each bundle declares a sequence of  
71 computations (its behavior) which will be executed in an **environment**. A bundle can *depend* on other  
72 bundles, and when being executed by an environment, its DAG of dependencies is linearized and run in  
73 sequence. In our use, bundles correspond to what an EC2 instance should do and an environment to an  
74 image (AMI: **A**maz**o**n **M**ach**i**ne **I**mag**e**) which prepares the basic configuration, downloads the Scala code  
75 and runs it.

## 76 2.4 2.x Datasets

77 Datasets is a Scala library developed by **so and so** to declare datasets and their locations. **Data** is  
78 represented as type-indexed fields: Keys are modeled as singleton types, and values correspond to what  
79 could be called a denotation of the key: a value of type `Location` tagged with the key type. Then a  
80 **Dataset** is essentially a collection of data, which are guaranteed statically to be different through type-level  
81 predicates, making use of the value type correspondence which can be established through singleton types  
82 and implicits. A dataset location is then just a list of locations formed by locations of each data member of  
83 that dataset.

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

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

## 93 2.5 2.x Loquat

94 Loquat is a library developed by **so and so** designed for the execution of embarrassingly parallel tasks  
95 using S3, SQS and EC2.

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

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

102 Both worker and manager instances are Statika bundles. In the case of the worker, it can declare any  
103 dependencies needed to perform its task: other tools, libraries, or data.

104 All configuration such as the number of workers or the instance types is declared statically, the  
105 specification of a loquat being ultimately a Scala object. There are deploy and resource management  
106 methods, making it easy to use an existing loquat either as a library or from (for example) a Scala REPL.

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

## 2.6 2.x Type-safe DSLs for BLAST and FLASH

We developed our own type-safe DSLs (Domain Specific Language) for FLASH and BLAST expressions and their execution.

### 2.6.1 2.x.a BLAST DSL

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

- all required arguments are provided
- only valid options are provided
- correct types for each option value
- valid output record specification

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

### 2.6.2 2.x.b FLASH DSL

In the same spirit as for BLAST,

## 2.7 2.x Bio4j

Bio4j is a data platform integrating data from different resources such as UniProt or GO in a graph data paradigm. We use the module containing the NCBI Taxonomy, and the use their Java API from Scala in the assignment phase.

# 3. RESULTS

## 3.1 Overview

To tackle the challenges posed by metagenomics big data analysis outlined in the Introduction,

- AWS resources in Scala (??) - A new approach to data analysis specification, management and specification based on working with it in exactly the same way as for a software project, together with the extensive use of compile-time structures and checks.
- Parallelization and distributed analysis based on AWS, with on-demand infrastructure as the basic paradigm - fully automated processes, data and cloud resources management.
- Static reproducible specification of dependencies and behavior of the different components using *Statika* and *Datasets* - Definition of complex pipelines using *Loquat* a composable system for scaling/parallelizing stateless computations especially designed for Amazon Web Services (AWS) - Modeling of the taxonomy tree using the new paradigm of graph databases (Bio4j). It facilitates

the taxonomic assignment tasks and the calculation of the taxa abundance values considering the hierarchical structure of taxonomy tree (cumulative values). - per-read assignment (??)

## 3.2 3.x 16S Reference Database Construction

Our 16S Reference Database is a curated subset of sequences from NCBI nucleotide database **nt**. This subset of 16S sequences was selected by similarity with the bacterial and archaeal reference sequences downloaded from RDP database (?). RDP unaligned sequences were used to capture new 16S sequences from **nt** using BLAST similarity strategies and, then, performing additional curation steps to remove sequences with poor taxonomic assignments to taxonomic nodes close to the root of the taxonomic tree. All the nucleotide sequences included in **nt** database has a taxonomic assignment provided by the genbank sequence submitter. NCBI provides a table (available at <ftp://ftp.ncbi.nlm.nih.gov/pub/taxonomy/>) to do the mapping of any Genbank Identifier (GI) to its Taxonomy Identifier (TaxID). Thus, we are based on a submitter-maintained taxonomic annotation system for reference sequences that supposes a sustainable system able to face the expected number of reference sequences that will populate the public global nucleotide databases in the near future. Another advantageous point is that we are based on NCBI taxonomy, the de facto standard taxonomic classification for biomolecular data (?). NCBI taxonomy is, without any doubt, the most used taxonomy all over the world and the most similar to the official taxonomies of each specific field. This is a crucial point because all the type-culture and tissue databanks follow this official taxonomical classification and, in addition, all the knowledge accumulated is referred to this taxonomy. In addition NCBI provides a direct connection between taxonomical formal names and the physical specimens that serve as exemplars for the species (?).

If metagenomics results are easily integrated with the theoretical and experimental knowledge of each specific area, the impact of metagenomics will be higher than if metagenomics progress in a disconnected research branch. This strategy for building our database allows substituting the 16S database by any other subset of **nt**, even by the complete **nt** database if it would be needed, for example, for analyzing shotgun metagenomics data.

## 3.3 3.x Bio4j and Graph Databases

## 3.4 3.x MG7 Pipeline Description

## 3.5 3.x Taxonomic Assignment Algorithms

### 3.5.1 3.x.y Lowest Common Ancestor based Taxonomic Assignment

For each read:

1. Select only one BLASTN alignment (HSP) per reference sequence (the HSP with lowest e value)
2. Filter all the HSPs with bitscore below a defined BLASTN bitscore threshold  $s_0$
3. Find the best bitscore value  $S$  in the set of BLASTN HSPs corresponding to hits of that read
4. Filter all the alignments with bitscore below  $p * S$  (where  $p$  is a fixed by the user coefficient to define the bitscore required, e.g. if  $p=0.9$  and  $S=700$  the required bitscore threshold would be 630)
5. Select all the taxonomic nodes to which map the reference sequences involved in the selected HSPs:
  - If all the selected taxonomic nodes forms a line in the taxonomy tree (are located in a not branched lineage to the tree root) we should choose the most specific taxID as the final assignment for that read
  - If not, we should search for the (sensu stricto) Lowest Common Ancestor (LCA) of all the selected taxonomic nodes (See Figure X)

In this approach the value used for evaluating the similarity is the bitscore that is a value that increases when similarity is higher and depends a lot on the length of the HSP

### 3.5.2 3.x.z Best BLAST hit taxonomic assignment

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

## 3.6 3.x Other

General approach. An analysis is defined as a software project. It can evolve in the same way. We can run the analysis in a test phase, review configuration and changes, etc. Key advantages of this approach are

- **Reproducibility** the same analysis can be run again with exactly the same configuration in a trivial way.
- **Versioning** The analysis is a software project so it goes through the same stages, there can be different versions, stable releases, etc.
- **Reuse** we can build standard configurations on top of this and reuse them for subsequent data analysis.
- **Decoupling** We can start working on the analysis specification, without any need for data in a much easier way.
- **Expresiveness and safety** choose only from valid Illumina read types, build default FLASH command based on that, ...

### 3.7 3.x Using MG7 with some example data-sets

We selected the datasets described in [Kennedy-2014] (??)

## 3.8 3.7 MG7 availability

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

# 4 4. DISCUSSION

## 4.1 4.1 Novelty points of MG7

The most innovative ideas and developments integrated in MG7 are:

- The management dependencies checking their correctness before compilation using Scala type system
- The automation of cloud resources and processes (parallelization management) - The cloud-oriented development of the system including a modeling AWS resources based on the powerful data typing of Scala
- The use of the Graph databases paradigm to store and manage the taxonomy tree to obtain the taxonomic assignments and the cumulative frequencies - MG7 provides a sustainable model for updating the database of reference sequences appropriate to face the challenging amount of sequences that are generating the new high throughput technologies of sequencing

## 4.2 4.2 Designed for future challenges

Other possible uses of the general schema: statika, loquat, ...

218 **4.3 MG7 Future developments**

219 4.3.1 4.3.1 Comparison of groups of samples

220 4.3.2 4.3.2 Interactive visualizations using the output files of MG7 (Biographika project)

**5 5 ACKNOWLEDGEMENTS**

221 INTERCROSSING (Grant 289974)

**6 6 REFERENCES**

**REFERENCES**

222 Zeller, G., Tap, J., Voigt, A. Y., Sunagawa, S., Kultima, J. R., Costea, P. I., et al. (2014). Potential of fecal  
223 microbiota for early-stage detection of colorectal cancer. *Molecular systems biology* 10, 766