# MG7: A fast horizontally scalable tool based on cloud computing and graph databases for microbial community profiling

Evdokim Kovach, Alexey Alekhin, Marina Manrique,
Pablo Pareja-Tobes, Eduardo Pareja, Raquel Tobes and
Eduardo Pareja-Tobes

Oh no sequences! Research Group, Era7 bioinformatics
Granada, Spain

April 8, 2014

Era7 bioinformatics

# What is metagenomics?

Metagenomics is the the study about collecting of genetic material from mixed community of organisms (usually bacteria):

- Soil samples
- Marine samples
- Clinical samples.

The typical problem of metagenomics is to obtain information about species composition in the sample:

- Which species are presented in the sample?
- How many different species are presented in the sample?
- How many species from the given genus are presented in the sample?

Era7 bioinformatics

16S rRNA gene is widely used to identify bacteria. There are several publicly available databases:

- NCBI 16S
- Greengenes
- SILVA.

MG7 is a cloud tool that performs assignment to the taxonomy tree based on mapping of reads from samples against 16S database.

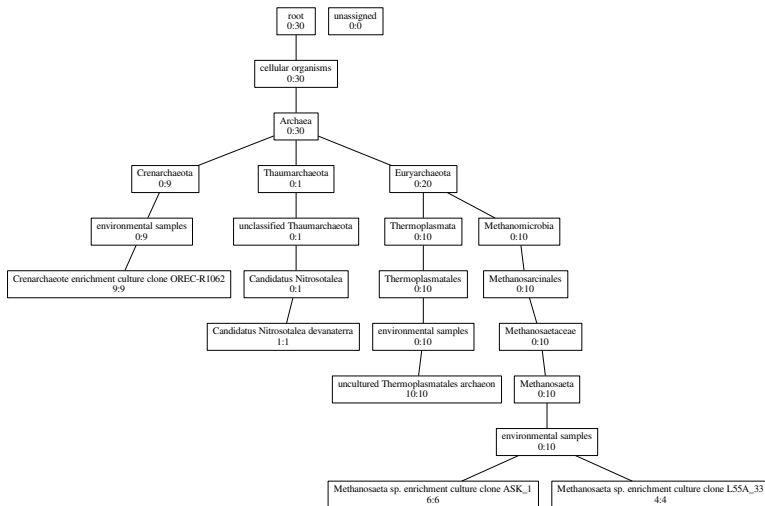Metapasta is a version of MG7 written in Scala that designed to process a really big amounts of metagenomics data.

Era7 bioinformatics

- Sequenced metagenomics samples – reads in FASTQ format
- Taxonomy database (we are using NCBI taxonomy database that is integrated to Bio4j)
- 16S database (by default Metapasta uses NCBI 16S with some filtering).

Era7 bioinformatics

As result Metapasta produces

- For every taxonomy id and sample:
  - number of reads from the sample that were assigned to this tax id
  - number of reads from the sample that were assigned to all successors of this tax id
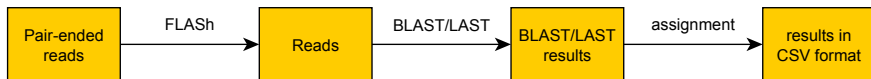- Aggregated data by samples.

# Results. Trees

# Results. Tables

| taxId | name | rank | supermock2.count | supermock2.acc | total.count | total.acc |
|---|---|---|---|---|---|---|
| total | | | 20 | 234 | 20 | 234 |
| 159447 | uncultured Corynebacterium sp. | species | 6 | 6 | 6 | 6 |
| 404941 | Mycobacterium salmoniphilum | species | 3 | 3 | 3 | 3 |
| 37637 | Corynebacterium pseudodiphtheriticum | species | 2 | 2 | 2 | 2 |
| 1221985 | Mycobacterium sp. ITM090653 | species | 2 | 2 | 2 | 2 |
| 319705 | Mycobacterium abscessus subsp. bolletii | subspecies | 2 | 2 | 2 | 2 |
| 1079047 | Mycobacterium sp. R5 | species | 1 | 1 | 1 | 1 |
| 43769 | Corynebacterium propinquum | species | 1 | 1 | 1 | 1 |
| 592914 | Corynebacterium sp. M71_S35 | species | 1 | 1 | 1 | 1 |
| 948102 | Mycobacterium franklinii | species | 1 | 1 | 1 | 1 |
| 1774 | Mycobacterium chelonae | species | 1 | 1 | 1 | 1 |
| 2 | Bacteria | superkingdom | 0 | 20 | 0 | 20 |
| 2037 | Actinomycetales | order | 0 | 20 | 0 | 20 |
| 131567 | cellular organisms | no rank | 0 | 20 | 0 | 20 |
| 1 | root | no rank | 0 | 20 | 0 | 20 |
| 85007 | Corynebacterineae | suborder | 0 | 20 | 0 | 20 |
| 1760 | Actinobacteria | class | 0 | 20 | 0 | 20 |
| 201174 | Actinobacteria | phylum | 0 | 20 | 0 | 20 |

Era7 bioinformatics

# Metapasta pipeline



- (optional) FLASh merging paired-end reads into bigger reads
- Mapping reads against the 16S database (with BLAST or LAST)
- Assignment to the taxonomy tree using Bio4j.

# Mapping problem

Mapping NGS reads against the 16S database requires really a lot of computational resources.
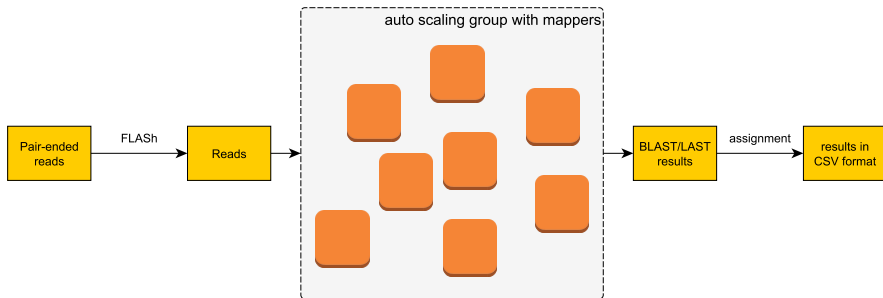
For example even on fast computers with SSD and big size of RAM mapping of one read with BLAST takes more than 0.2 seconds.

$$1000000 \times 0.2s \approx 56h$$

The mapping time can be improved by using a more efficient mapping tool (by default Metapasta uses LAST that in 100 times faster than BLAST in this case).

Metapasta uses AWS (Amazon Web Services) to perform all computations (EC2 instances):
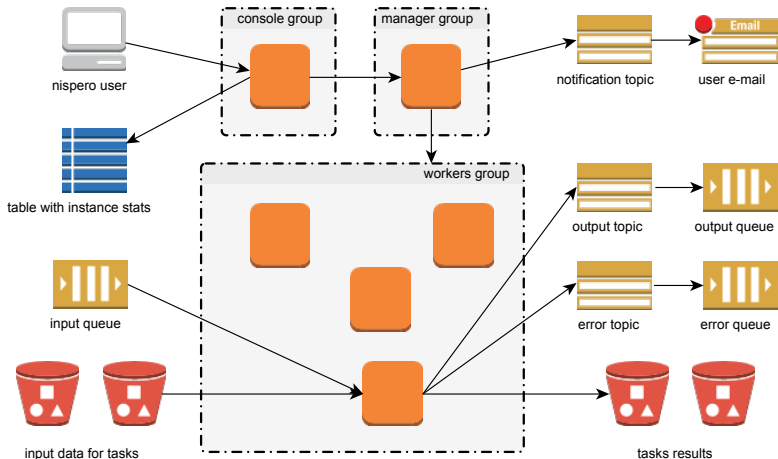
Besides computations Metapasta uses AWS for all data management:

- Reads from samples (S3)
- Taxonomy assignment tables and trees in PDF (S3)
- Metapasta can upload all reads (with assignments) to DynamoDB table.

Era7 bioinformatics

- Scala library for building distributed using AWS
- Easy to use (only AWS account is needed)
- Scalable and robust.

# Nispero. Architecture



nispero user

console group

manager group

notification topic

user e-mail

table with instance stats

workers group

input queue

output topic

output queue

error topic

error queue

input data for tasks

tasks results

Era7 bioinformatics

We are using idempotent commutative monoids to describe distributed systems in Nispero.

$$Reads \otimes Reads \xrightarrow{merge} Reads \xrightarrow{BLAST} AssignTable \otimes Reads$$

- They are powerful enough to build complex systems
- It is easy to work with them abstractly (for example to create graphical workflow editor that produces distributed systems).

Era7 bioinformatics

Do not lose opportunity to try this amazing Spanish fruit!

# Thank you for your attention!