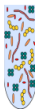
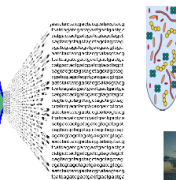
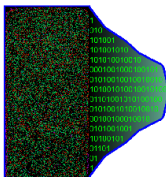


Introduction into the processing of raw data

Giuseppe D'Auria



FISABIO, Valencia



Norwich 12-16 October 2015

CONSIDERING NEEDED STORAGE SPACE BY TECHNOLOGY

Data Storage

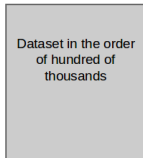
Size ranges

Sanger Sequencing



Datasets in the order of thousands of sequences

454



Dataset in the order of hundred of thousands

Illumina



Dataset in the order of millions of sequences

CONSIDERING NEEDED STORAGE SPACE BY TECHNOLOGY

Data Storage

Size ranges

Sanger Sequencing

454

Illumina

Solid

Dataset
order of t
of seq

Dataset in the order
of xxx of million of
sequences

CONSIDERING NEEDED STORAGE SPACE BY TECHNOLOGY

Data Storage

Size ranges

Sanger Sequencing

454

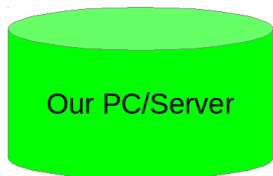
Illumina

New Illumina HiSeq systems

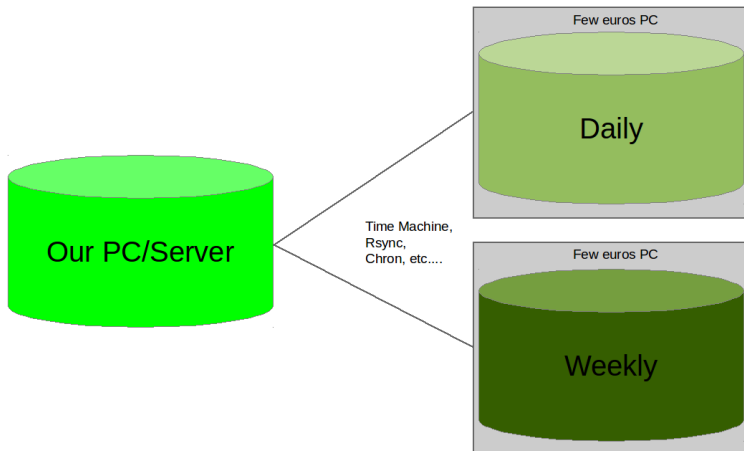
Dataset
order of t
of seq

Dataset in the order
of xxx of million of
sequences

CONSIDERING BACKUP



CONSIDERING BACKUP



CONSIDERING BACKUP

WIDE RANGE OF SOLUTIONS

BIG Servers, NAS systems, etc..



Thousands of Euro/Dollars

We spend so much money for sequencing, we can save few of them for saving our data

CONSIDERING BACKUP

WIDE RANGE OF SOLUTIONS

BIG Servers, NAS systems, etc..



Thousands of Euro/Dollars



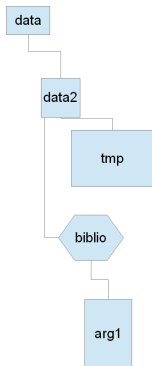
~ 50 Euros/Dollars
+ some tera of disks

We spend so much money for sequencing, we can save few of them for saving our data

DATA STORAGE AND DISK STRUCTURE



DATA STORAGE AND DISK STRUCTURE

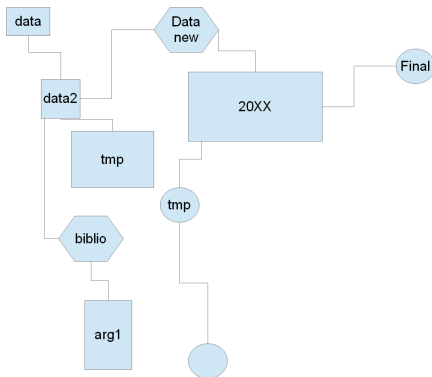


GENERALITAT
VALENCIANA

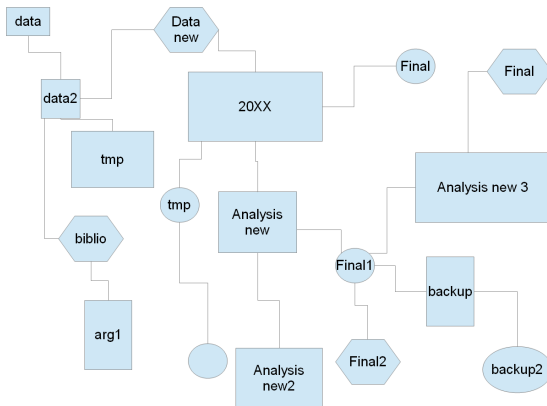


Fundación para el Fomento de la
Investigación Sanitaria y Biomédica
de la Comunitat Valenciana

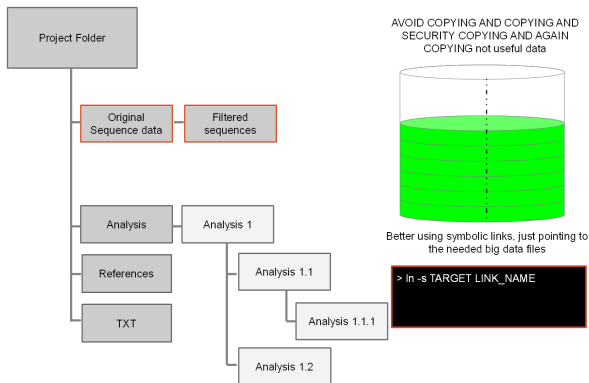
DATA STORAGE AND DISK STRUCTURE




DATA STORAGE AND DISK STRUCTURE



DATA STORAGE AND DISK STRUCTURE



THE SYSTEM

- Linux or Windows? 
- Both allow good bioinformatics analysis
- Linux is more stable for massive *data crunching* analysis and it is FREE
- Windows is not FREE
- Most of the software work in both systems but several are exclusively working on Linux.
- The best structure for bioinformatics (just my personal advice):
 - A Linux Desktop system (Ubuntu – Fedora) +
 - A virtual machine (Virtual Box)

DATA FORMAT

FASTA AND FASTAQUALITY FORMAT

FASTA

```
>G12OEMT03CWU1
AGAGTTTGATCATGGCTCAGGATGAACGCTAGCGGCAGGCCTAACACATGCAAGTCGAGGGAGGAG
CCTTCGGGCTTCGACCGCGTACGGGTGCGTAAAG
>G12OEMT03DH3XQ
AGAGTTTGATCATGGCTCAGTGCCAGCCGCCGCGGAGCGCATTAG
>G12OEMT03DD28C
AGAGTTTGATCCTGGCTCAGGGTGGTCATATGTTTGAATTGGTGCCAGCCGCCGCGGAGCGCATT
AG
>G12OEMT03DGC48
AGAGTTTGATCATGGCTCAGGAGGTGCCAGCAGCCGCGGAGCGCATTAG
>G12OEMT03C0MSF
AGAGTTTGATCCTGGCTCAGGACGAACGCTGGCGCGGTGCCTAATACATGCAAGTAGAACGCTGAA
GCTTGGCGCTTGCACCGAGCGGATG
```

DATA FORMAT

FASTA AND FASTAQUALITY FORMAT

FASTA

```
>G12OEMT03CWU1
AGAGTTTGATCATGGCTCAGGATGAACGCTAGCGGCAGGCCCTAACACATGCAAGTCGAGGGAGGAG
CCTTCGGGCTTCGACCGGCGTACGGGTGCGTAACTG
>G12OEMT03DH3XQ
AGAGTTTGATCATGGCTCAGTGCCAGCCGCCGCGGAGCGCATTAG
>G12OEMT03DD28C
AGAGTTTGATCCTGGCTCAGGGTGGTCATATGTTTGAATTGGTGCCAGCCGCCGCGGAGCGCATT
AG
>G12OEMT03DGG48
AGAGTTTGATCATGGCTCAGGAGGTGCCAGCAGCCGCGGAGCGCATTAG
>G12OEMT03C0MSF
AGAGTTTGATCCTGGCTCAGGACGAACGCTGGCGGCGTGCCTAATACATGCAAGTAGAACGCTGAA
GCTTGGCGCTTGCACCGAGCGGATG
```

QUALITY

```
>G12OEMT03CWU1
40 40 38 30 20 20 20 30 38 36 36 36 36 38 40 40 40 40 39 38 38 38 34 40 40 40 40 40 40
40 40 40 40 40 40 40 40 40 40 40 40 40 39 39 39 39 34 34 35 39 40 40 40 36 39 39 40 40 39 39
39 39 40 40 39 39 39 40 40 40 40 40 40 40 39 39 39 40 40 40 39 39 38 35 32 35 40 40 40 40
40 40
>G12OEMT03DH3XQ
40 40 40 38 20 20 20 30 38 36 36 36 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40
30 30 30 40 40 40 40 35 35 34 34 39 35
>G12OEMT03DD28C
40 38 37 35 22 22 22 26 31 35 36 33 30 32 33 36 36 30 28 20 18 18 35 27 30 32 32 32 27 21 22 16
16 14 19 19 23 23 23 23 23 21 24 27 32 27 27 25 27 30 24 24 25 27 26 28 28 32 22 29 27 25 22 20
19 21 27
>G12OEMT03DGG48
40 40 40 36 21 21 20 30 36 40 40 40 36 36 40 40 40 40 34 30 21 21 25 26 36 36 40 34 32 32 32
31 31 31 26 23 22 25 20 30 34 25 29 24 29 23 24
>G12OEMT03C0MSF
40 40 36 28 19 19 19 28 31 36 36 36 37 36 40 40 40 39 39 39 40 40 40 40 40 40 40 40 40 40
40 40 40 39 39 39 40 40 40 40 40 39 35 35 35 35 34 39 40 40 40 40 40 39 39 39 39 39
39 39 39 39 40 40 40 40 40 39 39 39 40 40 40 40 40 39 39 39 39
```


STANDARD FLOWGRAM FORMAT

GIUSEPPE D'AURIA

DATA FORMAT

FASTQ FORMAT

```

0AAII-22123:123:ABCDHFHT:4:1101:1985:2240 1:N:0:ATTTCT
ATCTGACGCCGCCGATTTGTATGCAGTAAATTTATTTATGAGCAGAGGCATA
+
000FFFBDDFBHBBHIIICBFHIIIGGIIIGGHHIGCHGHIIDHGIIIIIIGI
0AAII-22123:123:ABCDHFHT:4:1101:1969:2247 1:N:0:ATTTCT
TAAAGCCCGCCAGTTTGCATGCCAGGTCATGACAGAGGCATAAAACCGA
+
0C0FFFFFBHBBHJJJJJJJJJJJJJIFHHHIOIOIOIOIIIIIOIOJIEHH
0EAS139:136:FC706VIT:2:2104:15343:197393 1:Y:18:ATCAGS
GGAGTTTCATTACAAATTTATATATTTAAAGAGGNNNAGNNNNNINACTGAA
+
CC0FFFFFBHBBHPIIJJJJHHHIDHJJIFHHFBH03#1#1###00:DGFI
0AAII-22123:123:ABCDHFHT:4:1101:2226:2183 1:N:0:ATTTCT
TTCAGTTTGTGATGTGCGACGATGTTTGGCTCANGGCGCTNNNGTCTTGGG
+
CC0FFFBFBHBBHGGHGIIIJJJJJJJJJJJJJ0?B*-7#---;CHIJH
0AAII-22123:123:ABCDHFHT:4:1101:2094:2194 1:N:0:ATTTCT
CTCCACACATAACATAACGTTCCGCCAGTGGTATGCCACGNCNACGATGAGC
+
<?0DDDFHIBBDGCBGIIIDFCGDGCD?:D:0F?:?GHE#0?;CB00F
0AAII-22123:123:ABCDHFHT:4:1101:2544:2173 1:N:0:ATTTCT
GGCGCGGACGACGAAACACATCATCATGCTGCTCMNNAAGGAGGACGACGA

```

DATA FORMAT

FASTQ FORMAT

```

@AAII--Z2123:123:ABCDPEFHT;4:1101:1985:2240 :1:N:0:ATTCTC
ATCTGACGCCCGCATTTGTGCAGTAAATTATTTATATGACGACGAGGCATA
+
@@@FFBBDFFHHGGHIIICBFHIIIGGIIGGGHIGCHGHIDHGIIIIIIIIGI
@AAII--Z2123:123:ABCDPEFHT;4:1101:1969:2247 :1:N:0:ATTCTC
TAJACGCCCCGAGTTGCGATCCGAGSTGCATGACAGAGGCANTAJAACCGGA
+
GOCFFFFFHBBBJJJJJJJJJJJIFHHIJJIJJIJIIIIIIJJIJJJEHI
+
@EAS139:136:FC706VJ1;2:2104:15343:197393 :1Y:18:ATCAAG
GAGTTTCATTACAATTTATATATTTAAGAGGGBHBANGBNBRBGACTGAA
+
COO?FFFF?GHBUPIIJJJBHUIIUUFIKHFZB###1?####OO?DGFH
+
@AAII--Z2123:123:ABCDPEFHT;4:1101:2226:2193 :1:N:0:ATTCTC
TTCAGTTTGTGATGTCGCAOGATGGTTOGCTCANGOBCTNNNNTTCTCGG
+
COOFFEFFFHBBBHGGHGGIIJJJJJJJJJJJ#07B*-7##--:CHI JH
@AAII--Z2123:123:ABCDPEFHT;4:1101:2094:2194 :1:N:0:ATTCTC
CTCCACACTAACATAACOSTTCCCACAGSTGGATGCGCCAGNICACTAGAGC
+
<?0D?DDDFHHBBDGCBGLIDFCGDGC??D:C0F??:GHF#*07:CN#0F
@AAII--Z2123:123:ABCDPEFHT;4:1101:2544:2173 :1:N:0:ATTCTC
GGCGCCGACCTGAAAAACTCATCATGCTGATGCTANNNAITNNNNAGGACGA

```

SequenceID

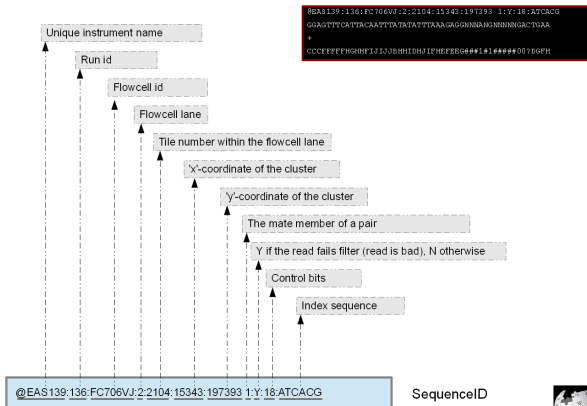
Sequence

- Optional

Quality

DATA FORMAT

FASTQ FORMAT - ID EXPLANATION



SequenceID



DATA FORMAT

FASTQ FORMAT - QUALITY SCORES

CCCCFFFFHGHFIJJJBHHIDHJIFHEFEEG####1#1#####00?DGFH

Quality

$$Q_{\text{phred}} = -10 \log_{10}(e)$$

e = estimated probability of a base being wrong

```

#####
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII
JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ
IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII
|'H$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNopqrstuvwxyz[\]^_`abdefghijklmnopqrstuvwxyz{|}~
33          59    64    73          104          126
0.....26...31.....40
      -5...0.....9.....40
      0.....9.....40
      3.....9.....40
0.....26...31.....41

S - Sanger      Phred+33,  raw reads typically (0, 40)
X - Solexa     Solexa+64, raw reads typically (-5, 40)
I - Illumina 1.3+ Phred+64, raw reads typically (0, 40)
J - Illumina 1.5+ Phred+64, raw reads typically (3, 40)
    with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (bold)
    (Note: See discussion above).
L - Illumina 1.8+ Phred+33,  raw reads typically (0, 41)

```



HAVE A LOOK AT THE PROJECT

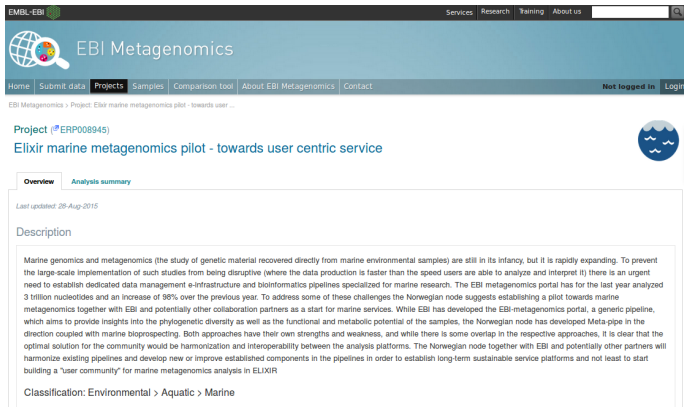
WHERE WE CAN FIND METAGENOMICS DATA

The screenshot shows the EBI Metagenomics website. At the top, there's a navigation bar with links: Home, Submit data, Projects, Samples, Comparison tool, About EBI Metagenomics, and Contact. Below this is a large banner with the text "Submit, analyse, visualize and compare your data." and a prominent "SUBMIT DATA" button. Under the banner, there are statistics: 7675 data sets, 4148 metagenomics, 390 metatranscriptomics, 3070 amplicons, and 67 assemblies. To the right, it shows 5578 runs, 4876 samples, and 136 projects. Further right, it displays 2097 runs, 2035 samples, and 91 projects. Below the statistics, there's a "Browse projects" section with a "By selected biomes" filter. This filter shows icons for Soil (17), Marine (28), Forest (4), Non-human host (33), Engineered (12), Freshwater (3), Grassland (3), Human gut (22), Air (1), and Wastewater (2). To the right of the biomes, there's a "Latest projects" section with three project highlights, each with a brief description and a "View more" link.

<https://www.ebi.ac.uk/metagenomics/>

HAVE A LOOK AT THE PROJECT

Just to start we will work on real metagenomics data downloaded from public data on EMBL-EBI Metagenomics



The screenshot shows the EBI Metagenomics website. The header includes the EMBL-EBI logo and navigation links: Services, Research, Training, About us. The main navigation bar has links: Home, Submit data, Projects, Samples, Comparison tool, About EBI Metagenomics, Contact. The page title is 'EBI Metagenomics'. Below the navigation bar, there's a sub-header 'Project (ERP008945)' and the project name 'Elixir marine metagenomics pilot - towards user centric service'. The page has tabs for 'Overview' and 'Analysis summary'. The 'Overview' tab is active. The description text reads: 'Marine genomics and metagenomics (the study of genetic material recovered directly from marine environmental samples) are still in its infancy, but it is rapidly expanding. To prevent the large-scale implementation of such studies from being disruptive (where the data production is faster than the speed users are able to analyze and interpret it) there is an urgent need to establish dedicated data management e-infrastructure and bioinformatics pipelines specialized for marine research. The EBI metagenomics portal has for the last year analyzed 3 trillion nucleotides and an increase of 96% over the previous year. To address some of these challenges the Norwegian node suggests establishing a pilot towards marine metagenomics together with EBI and potentially other collaboration partners as a start for marine services. While EBI has developed the EBI-metagenomics portal, a generic pipeline, which aims to provide insights into the phylogenetic diversity as well as the functional and metabolic potential of the samples, the Norwegian node has developed Meta-pipe in the direction coupled with marine bioprospecting. Both approaches have their own strengths and weakness, and while there is some overlap in the respective approaches, it is clear that the optimal solution for the community would be harmonization and interoperability between the analysis platforms. The Norwegian node together with EBI and potentially other partners will harmonize existing pipelines and develop new or improve established components in the pipelines in order to establish long-term sustainable service platforms and not least to start building a "user community" for marine metagenomics analysis in ELIXIR'. The classification is 'Environmental > Aquatic > Marine'.

<https://www.ebi.ac.uk/metagenomics/projects/ERP008945>

HAVE A LOOK AT THE PROJECT

Contact details

Institute: UIT
Name: Espen Robertsen
Email: Espen.m.robertsen@uit.no

Associated runs

Sample Name	Sample ID	Run ID	Experiment type	Version	Analysis results
Muddy	ERS624612	ERR695596	Metagenomic	1.0	Taxonomy Function ↓
		ERR695596	Metagenomic	2.0	Taxonomy Function ↓
Sandy	ERS624613	ERR695597	Metagenomic	1.0	Taxonomy Function ↓
		ERR695597	Metagenomic	2.0	Taxonomy Function ↓

Submitting Centre	Run Date	Platform	Model	Read Count	Base Count
University of Tromsø, NORWAY		ILLUMINA	Illumina MiSeq		
Library Layout PAIRED	Library Strategy WGS	Library Source METAGENOMIC	Library Selection size fractionation	Library Name unspecified	

Navigation

Read Files

This table contains the files for run ERR695596

[Download files](#)

Download: - of 1 results in [TEXT](#)

[Select columns](#)

Showing results 1 - 1 of 1 results

Study accession	Secondary study accession	Sample accession	Secondary sample accession	Experiment accession	Run accession	Tax ID	Scientific name	Instrument	Library layout	Fastq files (ftp)	Fastq files (galaxy)	Submitted files (ftp)	Submitted files (galaxy)	CRAM files (ftp)	CRAM files (galaxy)
PRJEB7944	ERP008945	SAMEA3168559	ERS624612	ERX640077	ERR695596	412755	marine sediment metagenome	Illumina MiSeq	PAIRED			Fastq file 1 Fastq file 2	Fastq file 1 Fastq file 2		

IERALITAT
ENCIANA

Fundación para el Fomento de la
Investigación Sanitaria y Biomédica
de la Comunitat Valenciana

CREATING PROJECT FOLDER STRUCTURE

OPERATIVE FOLDER

```
# go to practice folder.. linux is case sensitive  
cd Metagenomic  
  
# have a look at the folder  
ls -ltr  
  
# have a look at the tree  
tree
```

CREATING PROJECT FOLDER STRUCTURE

CREATING PROJECT FOLDER STRUCTURE

```
# Create project folder  
mkdir project  
  
# Go to project folder  
cd project  
  
# have a look, it should be empty.  
ls -ltr
```

CREATING PROJECT FOLDER STRUCTURE

CREATING PROJECT FOLDER STRUCTURE

```
# Create project folder
mkdir project

# Go to project folder
cd project

# have a look, it should be empty.
ls -ltr
```

LINKING ORIGINAL DATA FILE - 454 SFF FILE

```
# Make a symbolic link FROM - TO
ln -s ../original_data/*.fastq ./
```

which means:

make a symbolic links to all files from one folder up, original data, everything ending with “.fq.gz” in this folder.

CLEANING/TRIMMING ORIGINAL DATA

ILLUMINA FILES (AS WE EXPLAINED PREVIOUSLY) ARE USUALLY GZIP-PED

```
# just in case we need to unzip them before cleaning  
# gunzip file.fastq.gz
```

PRINSEQ-LITE FOR CLEANING DATASETS.

Always, always, always READ the Manual...

```
# to have an idea of the program.... read the manual  
prinseq-lite.pl -h  
  
# execute prinseq-lite.pl with a bunch of parameters  
prinseq-lite.pl -fastq sample_R1.fastq -fastq2 sample_R2.fastq -out_format 3 -out_good cleaned \  
-min_len 50 -trim_qual_right 20 -trim_qual_type mean -trim_qual_window 20 -out_bad null  
  
# have a look at the output  
ls -ltr
```

FASTQ STATISTICS BY FASTQ-STATS FROM EA-UTILS

FASTQ-STATS FROM EA-UTILS

Have a look at the input/output

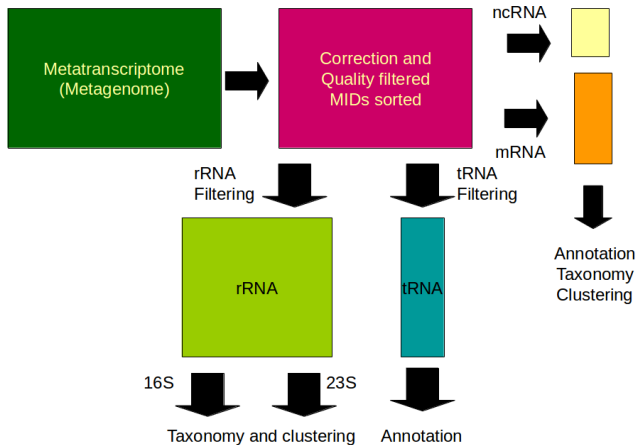
```
ls -ltr  
  
# have a look at the help  
fastq-stats -h  
  
# check fastq file after  
fastq-stats sample_R1.fastq  
  
# check fastq file after  
fastq-stats cleaned_1.fastq
```

JOIN PAIRED READS

FASTQ-JOIN FROM EA-UTILS


```
# so easy as reading the manual  
fastq-join -h  
  
# and calling it  
fastq-join -v ' ' cleaned_1.fastq cleaned_2.fastq -o sample.%.fastq  
  
# Have a look at the output  
ls -ltr
```

WE HAVE OUR DATASET!!!!



SEARCHING FOR RIBOSOMAL RNAs

We can use *SortMeRNA* which searches among provided databases.



SortMeRNA

[home](#) [help](#) [material](#) [FAQs](#)

DNA

- YASS
- Magnolia
- mreps
- ProCARs

HTS

- SortMeRNA
- CRAC
- Vidjil
- StoRM

RNA

- Carnac

SortMeRNA is a biological sequence analysis tool for filtering, mapping and OTU-picking NGS reads. The core algorithm is based on approximate seeds and allows for fast and sensitive analyses of nucleotide sequences. The main application of SortMeRNA is filtering rRNA from metatranscriptomic data. Additional applications include OTU-picking and taxonomy assignment available through QIIME v1.9+ (<http://qiime.org> - v1.9.0-rc1).

SortMeRNA takes as input a file of reads (fasta or fastq format) and one or multiple rRNA database file(s), and sorts apart rRNA and rejected reads into two files specified by the user. Optionally, it can provide high quality local alignments of rRNA reads against the rRNA database. SortMeRNA works with Illumina, 454, Ion Torrent and PacBio data, and can produce SAM and BLAST-like alignments.

If you use SortMeRNA, please cite:
Kopylova E., Noé L. and Touzet H., "SortMeRNA: Fast and accurate filtering of ribosomal RNAs in metatranscriptomic data", *Bioinformatics* (2012), doi: 10.1093/bioinformatics/bts611.

```
# we can format our database(s)
# indexdb_rna --ref ../rRNA_databases/silva-bac-16s-id90.fasta,./index/silva-bac-16s-db:\
# ../rRNA_databases/silva-bac-23s-id98.fasta,./index/silva-bac-23s-db
# I already did it

sortmerna --ref ../rRNA_databases/silva-bac-16s-id90.fasta,./index/silva-bac-16s-db:\
../rRNA_databases/silva-bac-23s-id98.fasta,./index/silva-bac-23s-db \
--reads sample.join.fastq --sam --num_alignments 1 --fastx --aligned sample_rRNA \
--other sample_not_rRNA --log -v
```


SEARCHING FOR tRNAs

From previously file **sample_non_rRNA** we can sort out tRNAs sequences using tRNAscanSE

Lowe Lab
tRNAscan-SE Search Server

Search for tRNA genes in genomic sequence



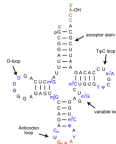
tRNAscan-SE 1.21

This web server is described in
[Schattner, P., Brooks, A.N., and Lowe, T.M. \(2002\) Nucleic Acids Res. 33, W686-689.](#)
The principles underlying the tRNAscan-SE program are described in
[Lowe, T.M. and Eddy, S.R. \(1997\) Nucleic Acids Res. 25, 953-964.](#)
If you use this tool in your investigations, please cite one of these references.

Instructions for using the tRNAscan-SE server and interpreting the output can be found in the [tRNAscan-SE README file](#).

If you would like to run tRNAscan-SE locally, you can get the UNIX [source code](#) (gzip'd tar file).

Analyzing tRNAs in a published genome? See our own tRNAscan-SE analyses of completed genomes in the [Genomic tRNA Database](#)



```
# first of all we need to convert fastq to fasta (FASTX-Toolkit)
fastq_to_fasta -Q 33 -i sample_not_rRNA.fastq > sample_not_rRNA.fasta
```

```
# executing tRNAscan-SE on general models (three kingdoms)
tRNAscan-SE -G -o tRNAs.txt sample_not_rRNA.fasta
```

```
# have a look at the results
less tRNAs.txt
```

```
# Extract the first column of the results (IDs column) skipping the first 4 lines
tail -n +4 tRNAs.txt | awk '{print $1}' > tRNAsIDs.txt
```

le la
dita
ono

SEARCHING FOR tRNAs

We have now to filter out **tRNAs.txt** reads in the IDs table from the **sample_non_rRNA.fastq**

THIS SCRIPT DOES NOT WORK!! SEE THE GOOD ONE IN
/HOME/TRAINING/BIN/REXTRACTFASTQFROMIDLIST.R

```
args<-commandArgs(TRUE)
if(args[1] == ""){
  print("usage: RExtractFastqFromIdList.R tabIdFile IN-FastqFile OUT-FastqFile");
  print("library ShortRead is required");
  print("Try again.....")
  q()
}
suppressMessages(library("ShortRead"))

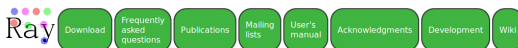
# Read table
ta<-read.table(args[1], sep="\t")

# Read fastq file
fq<-readFastq(args[2])
```

```
RExtractFastqFromIdList.R tRNAsIDs.txt sample_not_rRNA.fastq sample_no_trnas.fastq
```

ASSEMBLING METAGENOME

Now we are ready to assemble our filtered metagenome. We will use **Ray**



Ray -- Parallel genome assemblies for parallel DNA sequencing

```
# Execute SPAdes on sample_no_trnas.fastq
spades.py -s sample_no_trnas.fastq -o spades_out

tree -f

less .....
```



GENERALITAT
VALENCIANA



Fundación para el Fomento de la
Investigación Sanitaria y Biomédica
de la Comunitat Valenciana

SEARCHING FOR ORFs.

FOR ORFs SEARCH WE WILL USE **prodigal** SOFTWARE

We will search orfs within spades_out/contigs.fasta

```
prodigal -i spades_out/contigs.fasta -a orfs.faa -d orfs.fna -f gff -o orfs.gff -p meta
```