Microbiome Analysis

Introduction to sequencing technologies data pre-processing and quality control

La Serena 3-7.03.2025



Adam Ossowicki

Few things about myself



Traits involved in biological control of plant pathogens by *Pseudomoas*



IHSM

Microbiomes !!!



Sanger, <u>Illumina</u>, PacBio, 10X, Oxford nanopore....

- Homogenous samples
- Obligatory PCR step (specific)
- What do we use it for ?





https://emea.illumina.com/

NovaSeq 6000 Sequencing System :up to 6 Tb and 20B single reads in < 2 days



Cluster growth

Imaging

Sequencing

:.

TGCTAC

Base calling

TGCTAC

Base calling

Source: Illumina; Next Generation Sequencing Edited by Jerzy K Kulski, DOI: 10.5772/61657



What is Illumina? And how does it work?



		illaminar D SP	illemina:	Illuminar S2	Illeminer S4	
	Flow cell	SP	S1	S2	S4	
DNA	Lane number	2	2	2	4	
	M reads per lane	400	800	2'000	2'500	
	Total M reads	800	1'600	4'000	10'000	

DNA end polishing

Adaptor ligation

Size selection (optional) Amplification with

Input

barcode index

NGS (Illumina)

When do we barcode and why?

A flowcell

Source: Illumina; Next Generation Sequencing Edited by Jerzy K Kulski, DOI: 10.5772/61657

Why most sequencing we do is outsourced ?



https://emea.illumina.com/



Where is the sequencer ???





Nanopore vs. llumina

- Size of sequencers
- Longer reads
 - how much longer ?
- More sequencing errors ?

https://oxsci.org/pore-over-this-advances-in-dna-sequencing/

What is metagenomics?

Answer: Depends who you ask !

- Unspecific
- **Requires computing** infrastructure *
- Analysis requires (more) knowledge and experience
- More versatile



- Specific (is it ?)
- Not necessarily ٠ requires computing infrastructure *
- Analysis requires (less) knowledge and experience
- Gets deeper into microbiome *
- Less versatile

0

ITS

18S

taxonomy profiles,

relative abundance

* Most of the times

Gets deeper into microbiome... example



Let's sequence something...

What you give and what you get in return ?

b) a) 8% 16% Actinoallomurus Hydrogenispora 7% 14% Streptosporangiaceae Ruminococcus Nonomuraea Ruminiclostridium 6% 12% Streptomycetaceae Lachnospiraceae Streptomyces Clostridium sensu stricto 12 5% 10% Amycolatopsis Clostridia_UCG-014_ge Nocardioides Staphylococcus 4% 8% Verrucosispora Paenibacillus Micromonosporaceae Ammoniphilus 3% 6% Microbacteriaceae Planococcaceae Geodermatophilus Lysinibacillus 4% 2% Jatrophihabitans Bacillus Acidothermus unclassified Bacillaceae 2% Nocardia Aneurinibacillus Mycobacterium Tumebacillus Qiagen Zymo Qiagen Zymo

isolation

quality check 📫 shipping



https://doi.org/10.1007/s12223-021-00866-0

Effects of DNA preservation solution and DNA extraction methods on microbial community profiling of soil

Who is the 4th horseman ?

Source of the gel: JGI



Windows: WinMD5Free or certutil (command line) or other Linux/Mac: has in-built software called ... md5sum



Where (not) to store my data ?





- Backed-up storage servers
- Repositories
- Clouds ?

Computing power





To work with big data we need Computing infrastructure like servers and clusters.



What are they ?

Why do we need them ?

How do we use them ?



What are they ?

computer (n.) 1640s, "one who calculates, a reckoner, one whose occupation is to make arithmetical calculations," agent noun from <u>compute</u> (v.).

Lp		Nr ściezki					1988	Cytry		12	Nr ścieżki				1.4.4	Cylry	
	1	2		3	4	5	Litery	i znaki	Lp.	1	2	Π	3	4	5	Litery	i znaki
1	•	•	•				A	-	17	•	٠	•	•		•	Q	1
2	•		•		•	•	в	2	18		•	•		•		R	4
3		•	•		•		с	1	19	•		•	•			s	
4	•		•		•		D	kto tam	20			•			•	т	5
5	•		•				E	3	21	•	•	•	•			U	7
6	•		•	•	•		F	wolny	22		•	•	•	•	•	v	=
7		•	•		•	•	G	wolny	23	•	٠	•			•	w	2
8			•	٠		•	н	wolny	24	•		•	•	•	•	x	1
9			•				1	8	25			•			•	Y	6
ю	•	•	•		•		J	dzwonek	26	•		•			•	z	•
n		•	•	•	•		к	(27			•		•		powrot karetki	
12		•	•			•	L)	28		•	•				obrot	walka
13			•	•	•	•	м		29	•	•	•	•	•	•	lit	ery
14			•	•	•		N	,	30	•	•	•		•	•	cyfry i znaki	
15			•		٠	٠	0	9	31			•	•			odstęp	
16							P	0	32								

They are **just** computers....



www.etymonline.co https://www.blog-wajkomp.pl/polskie-komputery-odra-1103-odra-1204/

What are they ?

Laptop



performance

portability



easy to use







portability



easy to use



HPC – high performance computing







easy to use



Why do we need them ?

<u>Performance</u> <u>portability</u> or <u>easy to use</u>?

Why do they use Linux based operating systems*?





*most do but there are exceptions

Parameters !



cores



What I can do using...





What I can do using...

Server



Process big amplicon dataset

> Process and annotate moderate size metagenomes

Comparative genomics

What I can do using...



Do comparative metagenomics



How do we use them ?





Is it gonna work on my... ?

Windows







marcelm/**cutadapt**

Bow

- Most probably yes

- Probably no ...but you can emulate Linux in WSL (Windows Subsystem for Linux) and then probably yes

- Probably yes ... with exceptions



TIE

ne2





Windows: WinMD5Free or certutil (command line) or other Linux/Mac: has in-built software called ... md5sum



What you get in return...

- Sequence identifier, description
- Sequence
- **Optional field**
- Quality
- Next sequence

- @A00881:1079:HTJF2DRX2:1:210
 ATTGAGGAGTGTCAGCAGCCGCGGTAAT.

 - +

+

<u>न्यज्यज्यज्यज्यज्यज्यज्यज्यज्यज्य</u>ज्य<u></u> ज

Symbol	ASCII Code	Q-Score
6	54	21
7	55	22
8	56	23
9	57	24
	58	25
;	59	26
<	60	27
=	61	28
>	62	29
?	63	30
@	64	31
A	65	32
В	66	33
С	67	34
D	68	35
E	69	36
F	70	37
G	71	38
Н	72	39
L	73	40

@HWI-Mxxxx or @Mxxxx - MiSeq
@HWUSI - GAIIx
@HWI-Dxxxx - HiSeq 2000/2500
@Kxxxx - HiSeq 3000(?)/4000
@Nxxxx - NextSeq 500/550
@Axxxxx - NovaSeq
@Vxxxxx = NextSeq 2000
@AAxxxxx - NextSeq 2000 P1/P2/P3
@Hxxxxxx - NovaSeq S1/S2/S4

Phred quality score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1,000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%
60	1 in 1,000,000	99.9999%

Table 1- The probability values of a Phred quality score and the interpretation of their base-calling accuracy.

Quality is encoded Mostly using 2 different systems

Instrument identifier

run id flowcell id lane tiles ... Index sequence
@A00881:1079:HTJF2DRX2:1:2103:10348:4038 1:N:0:GCCAATAC+ACCGACAA
ATTGAGGAGTGTCAGCAGCCGCGGTAATACGGAGGGAGCTAGCGTTGTTCGGAATTACTGGGCGTAAAGCGC
+

Per base sequence quality



DATA PREPROCESSING - QUALITY CONTROL

Software: <u>FastQC</u> <u>FastP</u> FastX PRINSEQ TagCleaner

• • •

And very useful MultiQC for consolidation



DATA PREPROCESSING -ADAPTER TRIMMING, QUALITY FILTERING

 From QC data you may notice that adapters or primers (amplicon sequencing) are still present in your sequence.
 You should remove them either by providing the adapter/primer sequence or using a de-novo search.

- Recommended tools: Trimmomatic, Cutadapt, Fastp Dada2, bbduk, PRINSEQ++, AfterQC, Sickle, Flexbar, Seqtk...

Quality filtering – typical commands

• Trimmomatic

trimmomatic PE -phred33 raw_R1.fastq.gz raw_R2.fastq.gz trimmed_R1_paired.fastq.gz \
trimmed_R1_unpaired.fastq.gz trimmed_R2_paired.fastq.gz trimmed_R2_unpaired.fastq.gz \
ILLUMINACLIP:adapters.fa:2:30:10 LEADING:25 TRAILING:25 SLIDINGWINDOW:4:20 MINLEN:50
SLIDINGWINDOW:4:20 - Scans using a 4-base window, cutting when average quality < 20.</pre>

• Cutadapt

cutadapt -a <primer F> -A <primer R> -o trimmed_R1.fastq.gz -p trimmed_R2.fastq.gz \ -m 50 -q 25,25 -n 2 raw_R1.fastq.gz raw_R2.fastq.gz

• Fastp

fastp -i raw_R1.fastq.gz -I raw_R2.fastq.gz -o trimmed_R1.fastq.gz -O \
trimmed_R2.fastq.gz --detect_adapter_for_pe --cut_right --cut_mean_quality 25 \ -length_required 50 --n_base_limit 0 --html fastp_report.html --json fastp_report.json

examples

6th Plant Microbiome Symposium

microBl

#6

Málaga 2025

3-7 November 2025 Antequera, Málaga, Spain

Topics

Computational biology & Microbiomes

MicroS

0

Registration deadline: **April 2025**

Program outline:

- International keynote speakers
- Poster sessions
- Networking
- Joint dinners and excursions

Registration open now at:

6thplantmicrobiomesymposium2025.com



Thank you for your attention

