

DTU



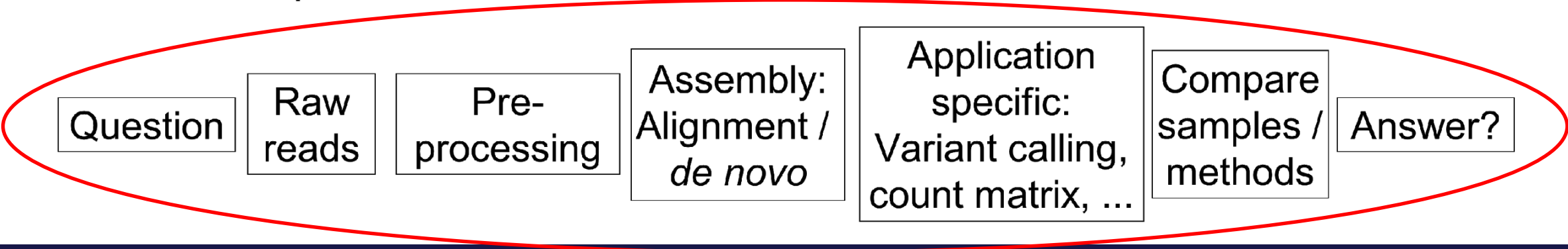
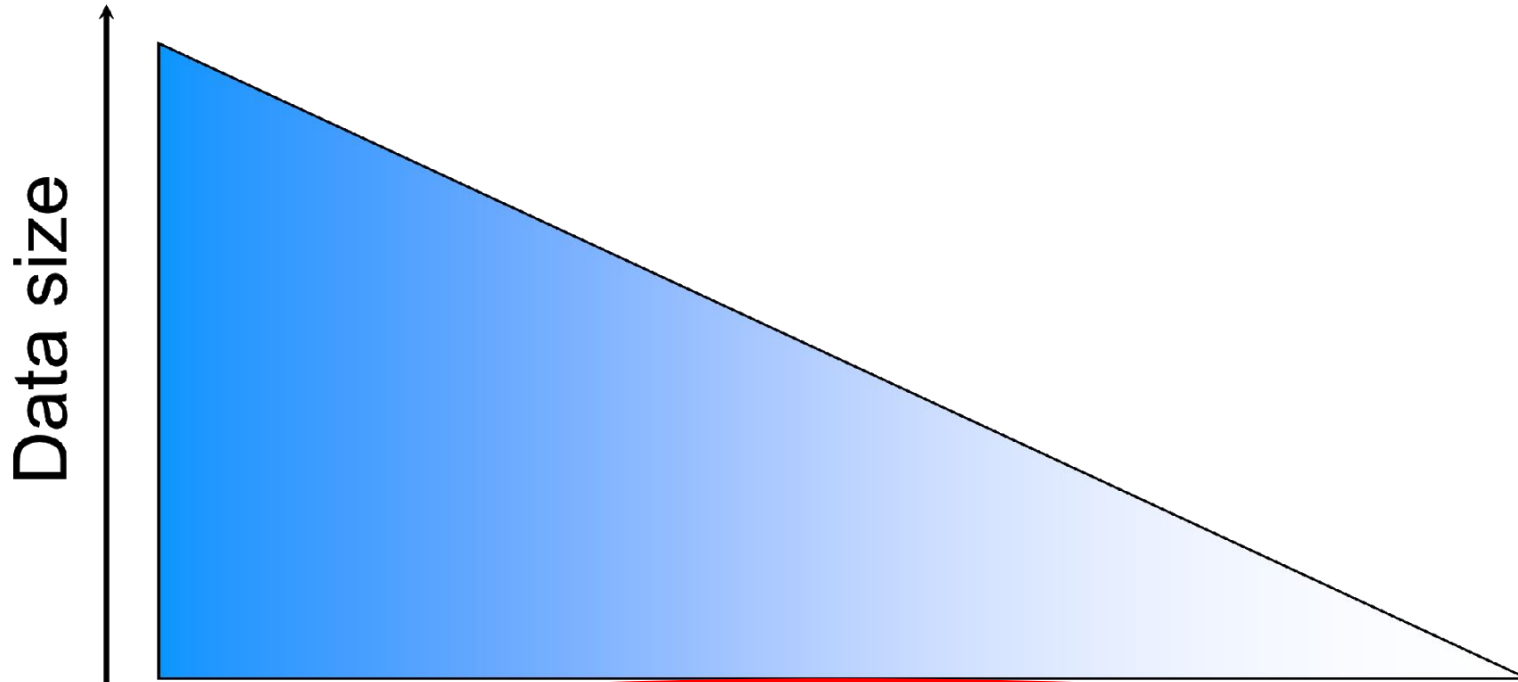


**DTU Health Technology
Bioinformatics**

Projects

*Gabriel Renaud
Associate Professor
Section of Bioinformatics
Technical University of Denmark
gabre@dtu.dk*

Generalized NGS analysis



Remember the slide from day 1? About the paragraph from a scientific paper?

Why are we here?

WES and WGS trio analysis. WGS sequencing and analysis for F01–08 and F13–20 were performed as described previously^{13,37}. Exome capture and sequencing of F09–12 were performed at the New York Genome Center (Agilent Human All Exon 50 Mb kit, Illumina HiSeq 2000, paired-end, 2 × 100) and the Broad Institute (Agilent Sure-Select Human All Exon v.2.0, 44-Mb baited target, Illumina HiSeq 2000, paired-end, 2 × 76). Sequencing reads were aligned to the hg19 reference genome using BWA (v.0.7.8). Duplicates were marked using Picard's MarkDuplicates (v.1.83, <http://broadinstitute.github.io/picard>) and reads were realigned around insertion/deletions (InDels) with GATK's IndelRealigner. Variant calling for SNVs and InDels was performed according to GATK's best practices by first calling variants in each sample with HaplotypeCaller and then jointly genotyping them across the entire cohort using CombineGVCFs and GenotypeGVCFs. Variants were annotated with SnpEff (v.4.2) and SnpSift (v.4.2), and allele frequencies from the 1000 Genomes Project and the Exome Aggregation Consortium (ExAC)³⁸. De novo variants were called for probands using Triodenovo (v.0.06) with a minimum de novo quality score of 2.0 and subjected to manual inspection. Variants from F01–F08 were further

Learning objectives

1. Are you able to:
 - a. work in group and delegate tasks?
 - b. set realistic objectives?
 - c. use the command line?
 - d. understand the strength and weakness of each tool?
 - e. explain key steps in a critical manner?

Projects

- Try to analyze an empirical dataset and present results on poster
- Aim for at least 1 figure, 1 table or 2 figures
- 4-5 pr. group
- You can find a dataset on SRA/ENA
- Try to find raw data, untrimmed
 - If not, please contact us

Projects

- You can use your own data if everyone in the group agrees ***and*** it can be presented on a poster
- Subset! Do not analyze very large datasets (time, resources)
- Subset! Do not replicate every figure/table!

Group formation

- Try to create groups with multiple competences
- Chose a group based on eg. field of interest

- Do not bite off more than you can chew:
 - Downloading the data, preprocessing, aligning will take several days

Previous projects

1. Introduction

Preterm babies are often administered early extended antibiotic therapy[1]. These therapies have potential detrimental effects on gut microbiota and on development of antibiotic resistance (AR) genes. It is therefore critical to understand the impact of such a therapy on the gut of a preterm infant. A 2016 study[2] investigated 401 stool samples from 84 preterm babies, taken during the first months of life. In this project, we analyse a subsample of this dataset in an attempt to find out how the administration of antibiotics affects the development of the gut microbiome in preterm infants.

2. Data specifications

A subsample of the full 401 samples was obtained by selecting 3 babies who had been treated with antibiotics (case) and 3 who had not (control). Six samples with similar sampling profile was chosen to minimize impacts from variables other than antibiotic treatment such as diet and gestational age at birth[2]. The resulting subset totalled approximately 6 Gbases from Illumina paired end reads.

4. Workflow

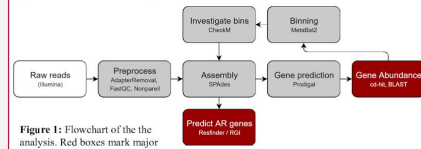


Figure 1: Flowchart of the analysis. Red boxes mark major result output.

3. Materials and methods

An initial run of FastQC was performed to evaluate the quality of the data (not shown), after which the reads were trimmed using the AdapterRemoval program. The coverage of the preprocessed genes was estimated using Nonpareil Curves (Figure 2).

Afterwards, the trimmed reads were assembled sample-wise using SPAdes, and the resulting contig files were analysed for resistance genes in ResFinder and in Resistance Gene Identifier (RGI) (Figure 3 & 4)

The contigs from the assembly were searched for bacterial genes using Prodigal and binned using MetaBat2. The binning result was analysed in CheckM (not shown), while the Prodigal output was used to create a species count matrix using cd-hit. Finally, the difference in species abundances between the samples were plotted (Figure 5). For a visual overview of the workflow see the flowchart (Figure 1).

5. Full coverage in samples

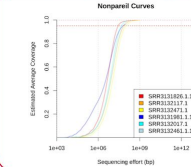


Figure 2: Using the Nonpareil curves we are able to estimate full coverage for all six samples. Furthermore, since the curves are closely grouped, the difference in diversity is estimated to be little.

6. Difference in resistome (RGI)

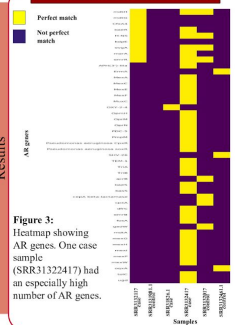


Figure 3: Heatmap showing AR genes. One case sample (SRR3132417) had an especially high number of AR genes.

7. Resistomes (ResFinder)

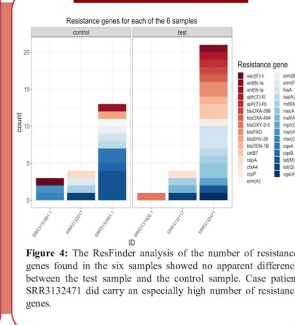


Figure 4: The ResFinder analysis of the number of resistance genes found in the six samples showed no apparent difference between the test sample and the control sample. Case patient SRR3132417 did carry an especially high number of resistance genes.

8. Varying microbiomes

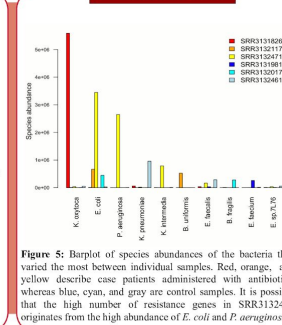


Figure 5: Barplot of species abundances of the bacteria that varied the most between individual samples. Red, orange, and yellow describe case patients administered with antibiotics, whereas blue, cyan, and gray are control samples. It is possible that the high number of resistance genes in SRR3132417 originates from the high abundance of *E. coli* and *P. aeruginosa*.

9. Abundant bacteria of interest

Sample nr.	Bacteria	AR resistance	Administered Antibiotic(s)	Potential diseases
SRR3131626.1	<i>Klebsiella oxytoca</i>	OXY-2-4 (beta-lactam)	Ampicillin (type of beta-lactam)	"Bronchopneumonia, urinary tract infections and septicemia" [3]
SRR3132461.1	<i>Klebsiella pneumoniae</i>	oxyA, oxyB (beta-lactam)	Control (antibiotic treatment at birth only)	"Nosocomial and systemic infections" [5]
SRR3132471	<i>P. aeruginosa</i>	SxvR, CspR (coding for resistance to 19 classes of antibiotics)	Vancomycin, Gentamicin, Meropenem, Colistin	"Urinary tract infections, ventilator associated pneumonia and infections related to mechanical heart valves, shunts, grafts and stents" [6]

Table 1: A selection of three of the bacteria which were found in high abundance (Figure 5). Two of these bacteria have resistance to the administered antibiotics.

10. Conclusions & Future perspective

- Analysis of our assembly using MetaBat2 and CheckM resulted in large and non-specific bins. This could indicate an error in our assembly, but due to time limits we were unable to redo this step.
- Investigation of the resistome using ResFinder and RGI identified a high number of AR genes in both case and control samples, with one case sample having more AR genes than the other. However, we did not attempt to prove statistically if the number of AR genes and antibiotic treatment are correlated.
- Identification of variation in species abundance between samples, determined using Prodigal and cd-hit, revealed that two case samples had an increased abundance of bacteria unique to those samples that have implications in disease.
- Perspective:* The pipeline shows promise, however, we were unable to draw any significant conclusion from our limited dataset. The gut microbiome of preterm babies is influenced by factors such as diet and gestational age[2]. Even though our subsample was selected with this in mind, prevalent high variability between samples persisted and a larger sample size is most likely needed in order to reveal how antibiotics modulate the gut microbiome and resistome of preterm infants.

References

[1] Tran RK, Bloom BT, Sigurd AR, Gerstmann DR (2016). Reported medication use in the neonatal intensive care unit: data from a large national data set. *Pediatrics*, 117: 1979-1987
 [2] Gilson, M, K. Wang, B. Ahnadi, S. Burroughs, C. A. Tarr, P. I. Vanner, B. B. & Durkin, G. (2016). Developmental dynamics of the preterm infant gut microbiota and antibiotic resistance. *Nature microbiology*, 1: 16024.
 [3] Singh, L, Carrique-Mun, J. & Nishi, M. (2016). Nosocomial infection: An emerging global health threat. *Medical Journal Australasian Federation*, 152: 289-295.
 [4] Nordmann, P, Cloutier, G. & Nasse, T. (2009). The real threat of *Klebsiella pneumoniae* carbapenemase-producing bacteria. *The Lancet Infectious diseases*, 9(4), 228-236.
 [5] Chari, B, K., Sazonov, M., Wertz, J. E., Kortright, K. E., Narayan, D., & Turchi, P. E. (2016). Plasmid selection restores antibiotic sensitivity in MDR *Pseudomonas aeruginosa*. *Scientific reports*, 6: 26717.

Posters

- Each group will create a poster

- ~~• You can print posters at the DTU library for 20-30kr~~

Online this year: send us a high resolution PDF!

Posters

1) The group number, student names and student numbers of all group members, must be stated on the poster

2) The poster must specify the individual students contribution to the project. It is allowed to state that everyone contributed equally

~~3) The poster must not extend the poster board (160 cm high, 120 cm wide)
Note, If you print through us the poster dimensions are: 1189mm x 841mm~~

[4\) Guide for making an good poster](#)

<http://wiki.bio.dtu.dk/teaching/index.php/Poster>

Grouping & Guidance

- Fill in group information in Google doc
- 5 min presentation tomorrow at 13
 - What do you plan to do?
 - How much data?
- Project assistance: every day
 - Teachers+TA via Discord
- Data goes here:

`/data/shared/groups/group_X`

Be nice!

- Run larger programs on the servers using nice eg.

```
nice -n 19 blastall -i alldatainthegalaxy -db everythingeversequenced
```

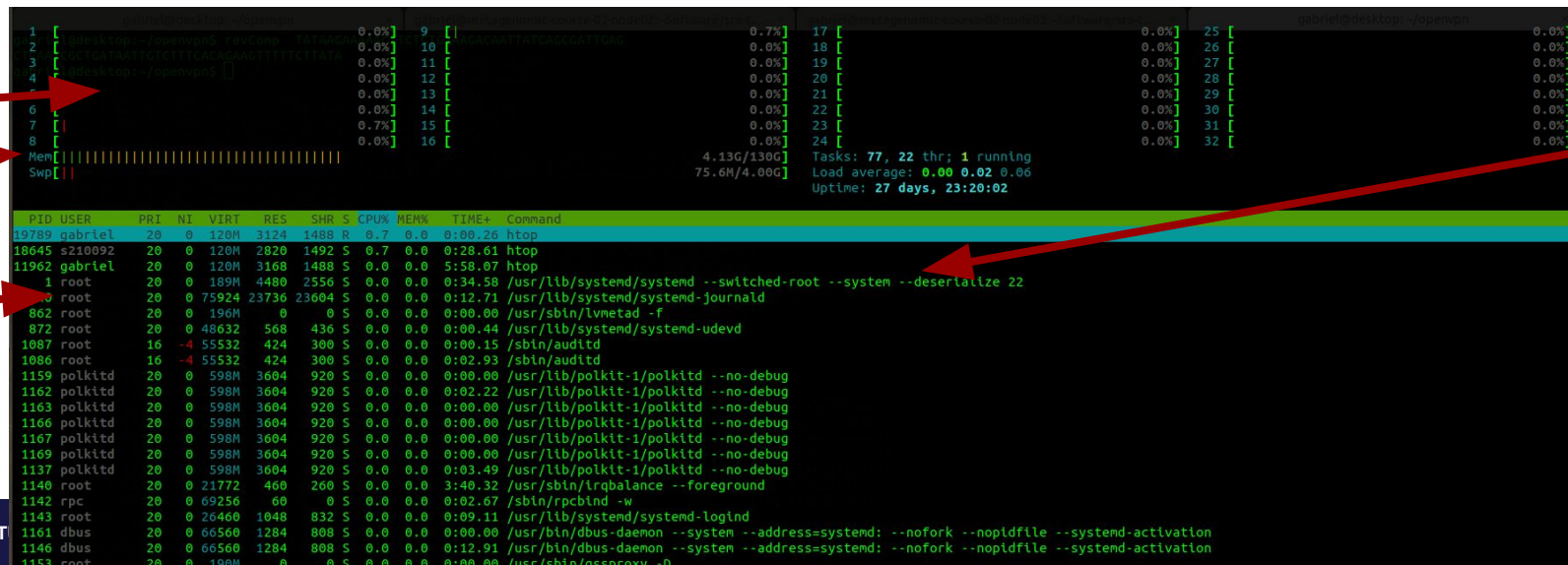
- How much memory am I using?

htop

CPUs

MEM

users



processes

Thou shall keep your files zipped

- Zip your vcf, text whatever files
 - there are tools to work with zipped files (zcat, zgrep, zless)
- Use BAM/CRAM never sam
- Beware, what is wrong with this?:

```
bwa mem reference.fasta input.fastq.gz > output.bam
```

Evaluation: presentation and oral exam

- You will give a group presentation about your poster (5-7 minutes)
 - each person should speak at least once.
 - what the study was about
 - what you have done
 - results you got:
 - Quality of data, replicate certain results, pitfalls
 - Please turn on your camera, we cannot evaluate you otherwise

Evaluation: presentation and oral exam

- We will ask once person at a time to come and we will ask you about 4-5 questions about the project:
 - The goal: did you understand what we taught in class and what you did
 - We can quiz you on your project and can have notions of what we saw in class
 - Do not memorize, understand!
 - Do not communicate with others in your group
- 2 evaluators will meet and the final mark will be a blend of your oral exam, group performance (minor tech talks) can help us distinguish between a 10 or 12.

Parting words

- No one size fits all solution for everything
 - How to genotype, population geneticists vs medical field
- Every tool shown in this class may/will be outdated in 5 years
 - Sorry for no textbook but it would be outdated soon
 - Read recent papers, reviews
 - bioRxiv is great but not peer-reviewed
- Question existing methods, pipelines, be wary of:
 - “This is how we do things around here”
 - “This is the standard pipeline for this kind of data”
- Understand how tools work, test