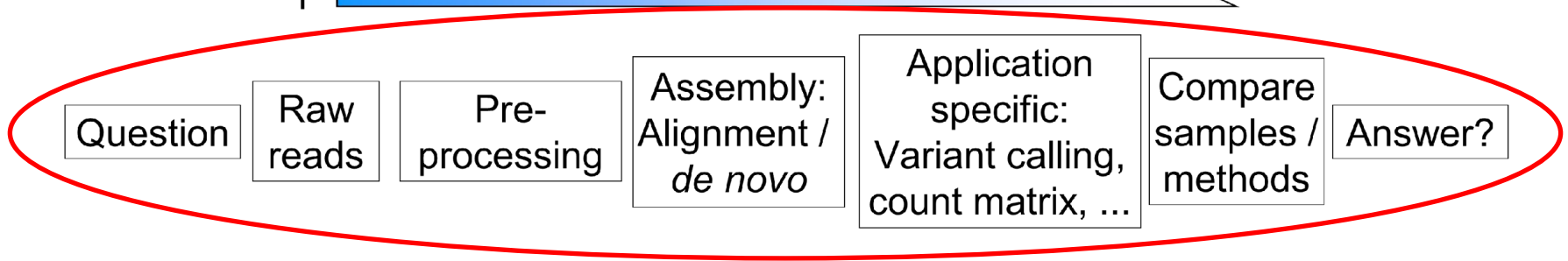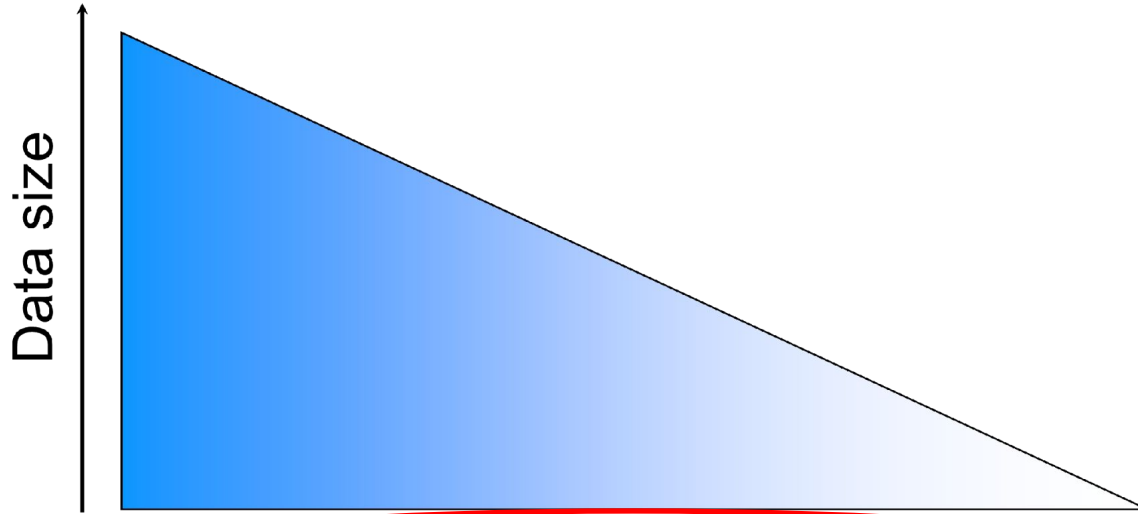**DTU Health Technology**
**Bioinformatics**

# *Projects*

*Gabriel Renaud*
*Associate Professor*
*Section of Bioinformatics*
*Technical University of Denmark*
*gabriel.reno@gmail.com*

# Generalized NGS analysis



Data size

Question | Raw reads | Pre-processing | Assembly: Alignment / *de novo* | Application specific: Variant calling, count matrix, ... | Compare samples / methods | Answer?

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

# Why are we here?

For single nucleotide variant calling, the data processing pipeline for detecting variants in Illumina HiSeq data is as follows. First the FASTQ files are processed to remove any adapter sequences at the end of the reads using cutadapt (v1.6). The files are then mapped using the BWA mapper (bwa mem v0.7.12). After mapping the SAM files are sorted and read group tags are added using the PICARD tools. After sorting in coordinate order, the BAMs are processed with PICARD MarkDuplicates. The marked BAM files are then processed using the GATK toolkit (v 3.2) according to the best practices for tumour normal pairs. They are first realigned using ABRA (v 0.92) and then the base quality values are recalibrated with the BaseQRecalibrator. Somatic variants are then called in the processed BAMs using muTect (v1.1.7) for single nucleotide variant and the Haplotype caller from GATK with a custom post-processing script to call somatic indels.

"In theory, there is no difference between theory and practice
but in practice there is a huge difference between theory and practice."

-Yogi Berra?
-Richard Feyman?
Someone else?

# Problems with NGS…

1. Reanalyzing previous data is part of the job

2. Impossible/difficult to reproduce results

3. "in house scripts" do not get me started…

4. Getting the data in the first place

5. Getting the metadata to match the data

# 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

  - Either replicate some results or ask your own question

- Aim for at least 1 figure, 1 table or 2 figures **of results**

- 5-6 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!

# Pitfalls

- Beware of esoteric tools or overuse of "in-house" scripts or programs

- No, you do not need to use exactly the same tools they used

- Some data (esp. human patient data) is not available due to privacy

  issues

- You can download subsets using SRA

- No single-cell data

# Pitfalls

- Limit the number of samples

- You will **not** get penalized if you gave a *bona fide* attempt at

  replicating results and cannot

- You can ask your own questions!

# 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

# Group formation

- Everyone is expected to put in work

- If there are freeloaders, let me know

# Previous projects



15

# Posters

- Each group will create a 1 page poster (not a report)

- You have to send it to:
  - gabriel.reno@gmail.com

Before Thursday 23rd at noon (12PM)

It is 25% of the grade.

# Posters

- I cannot read long paragraphs, please use very

  concise prose

- Bullet points are great!

- **Important text in bold**

# Posters

Little text, bullet points are amazing!

# Grouping & Guidance

- Fill in group information in Google doc

- 3 min chat tomorrow at 13

  ○ What do you plan to do?

  ○ How much data?

- Project assistance: every day

  ○ Teachers+TA

  ○ Write on Discord

# Grouping & Guidance

- Data goes here:

```
/home/projects/22126_NGS/projects/groupX
```

- Put the parsing scripts+plots there.
- Do not fabricate results.

# Pupil servers

| name | CPUs | RAM |
|------|------|-----|
| pupil1 | 40 | 252G |
| pupil2 | 24 | 110G |
| pupil3 | 24 | 94G |
| pupil4 | 16 | 63G |
| pupil5 | 16 | 63G |

# Be nice

- Run larger programs on the servers using nice eg.

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

- If the server is too full come talk to me :-)

- You can now all nodes ! But check available CPU/RAM. How to check?

CPUs

MEM

processes

users

# 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/vcf.gz **never sam/vcf**

- Use symlinks instead of copying (no wasted space!):

  ```
  ln -s /path/to/largefile .
  ```

- Beware, what is wrong with this?:

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

**2023:**

**G** What is wrong with this command line:
bwa mem  reference.fasta  input.fastq.gz > output.bam

It is difficult to determine what is wrong with this command line without more information. However, a few things that could be causing issues include:

- The reference.fasta and input.fastq.gz files might not be in the same directory as the command, or might not exist at all.
- The output file, output.bam, might not be writable (e.g. if the user does not have permission to write to the directory where the output file is located).
- There might be an issue with the specific version of bwa being used.
- You may need to provide more options for the 'bwa mem' command or it would cause an error.

It's recommended to check the existence of the input files, if you have the permission to write at the directory that you are trying to write the output, and consult the bwa manual to be sure that you have provided the necessary options.

**2024:**

**ChatGPT**

The command line you provided for running BWA (Burrows-Wheeler Aligner) contains a small error related to the output format. Here's the command you provided:

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

The issue with this command is that it directs the output of `bwa mem` to a file named `output.bam`, but the output is in SAM format, not BAM format. The BWA `mem` command generates output in SAM (Sequence Alignment/Map) format, which is a text format. The `.bam` file extension is typically used for binary alignment/map format, which is a binary version of SAM.

To correct this, you have two options:

1. **Output in SAM Format:** If you want the output in SAM format, you should change the output file's extension to `.sam` to reflect the correct format:

```shell
bwa mem reference.fasta input.fastq.gz > output.sam
```

2. **Convert to BAM Format:** If you need the output in BAM format, you should pipe the SAM output of BWA into a tool like `samtools` to convert it to BAM format. Here's how you can do it:

```shell
bwa mem reference.fasta input.fastq.gz | samtools view -bS - > output.
```

In this command, `samtools view -bS -` ...ads the SAM format from standard input

# How to get help?

- You want software installed, ask me or TAs
- Help with NGS software:



read error message

Google/ChatGPT

peers/discord

me/TAs

# The final exam

- Friday, Jan 24th at 9am

- I should be around

- Multiple choice

- 75% of the grade

# The final exam

- Understand the material

- Focused on the core lectures

- 3-4 general questions on the specialized lectures

# Things to remember

- I am not big on memorizing but…

- … please remember

  - Names of sequencing techs: e.g. Illumina, ONT

  - Their idiosyncrasies

  - Some basic formats we saw (e.g. FASTA, BAM, BED…)

# Parting words

- Do not trust your data, use:
  - quality controls
  - visualizations

- 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, benchmarking studies
  - bioRxiv is great but not peer-reviewed

# Parting words

- 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

- Do not trust your code, test

# Parting words

- Know the history of the seq. tech you are using as you may work with older data

- Do your literature search, use existing tools when possible
  - Google Scholar
  - Twitter (yes unfortunately)

- Talk to lab techs, molecular biologists, try to learn as much as possible about the biology and biotech behind the data being generated.

# Special projects/Master's projects

- Like NGS? Genotyping? Population genetics? Ancient DNA? pangenome graphs? sedimentary DNA?



- me
- 1 Postdoc
- 2 PhDs
- 4 Master's

The Modern and Ancient Genomes Group

# Articles in 2024

Issues | More content ▼ | Submit ▼ | Alerts | About ▼ | Molecular Biology and Evo

JOURNAL ARTICLE

## soibean: High-Resolution Taxonomic Identification of Ancient Environmental DNA Using Mitochondrial Pangenome Graphs 🔓

Nicola Alexandra Vogel ✉, Joshua Daniel Rubin, Anders Gorm Pedersen, Peter Wad Sackett, Mikkel Winther Pedersen, Gabriel Renaud ✉    Author Notes

**Article Contents**

Abstract
Introduction
Results
Di
Me

📄 PDF | ▥ Split View | 66 Cite | 🔧 Permissions | ◁ Share ▾

**Abstract**

New Results                                    🔔 Follow this prepri

## SAFARI: Pangenome Alignment of Ancient DNA Using Purine/Pyrimidine Encodings

ⓘ Joshua Rubin, ⓘ Jan van Waaij, ⓘ Louis Kraft, ⓘ Jouni Sirén, ⓘ Peter Wad Sackett, ⓘ Gabriel Renaud
doi: https://doi.org/10.1101/2024.08.12.607489

💬 0 | ☑ 0 | 👥 0 | ⚙ 0 | 🖵 0 | ▤ 0

Abstract | Full Text | Info/History | Metrics | 🗋 Preview PDF

### Abstract

Aligning DNA sequences retrieved from fossils or other paleontological artifacts, referred to as ancient DNA, is particularly challenging due to the short sequence length and chemical damage which creates a specific pattern of substitution ($C \rightarrow T$ and $G \rightarrow A$) in addition to the

New Results                                    🔔 Follow this preprint

## CarpeDeam: A *De Novo* Metagenome Assembler for Heavily Damaged Ancient Datasets

ⓘ Louis Kraft, ⓘ Johannes Söding, ⓘ Martin Steinegger, Annika Jochheim, Peter Wad Sackett, ⓘ Antonio Fernandez-Guerra, ⓘ Gabriel Renaud
doi: https://doi.org/10.1101/2024.08.09.607291

💬 0 | ☑ 0 | 👥 0 | ⚙ 0 | 🖵 0 | ▤ 0

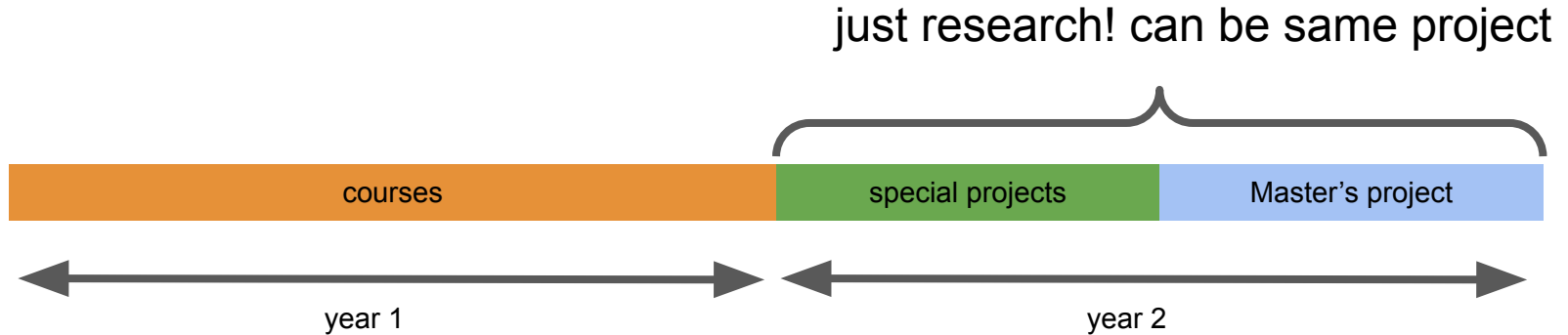Abstract | Full Text | Info/History | Metrics | 🗋 Preview PDF

### Abstract

*De novo* assembly of ancient metagenomic datasets is a challenging task. Ultra-short fragment size and characteristic postmortem damage patterns of sequenced ancient DNA molecules leave current tools ill-equipped for ideal assembly. We present CarpeDeam, a novel damage-aware *de novo* assembler designed specifically for ancient metagenomic samples. Utilizing

# Special projects/Master's projects

- Like research? Thinking about going to PhD? want to write a paper?

just research! can be same project

| courses | special projects | Master's project |
|---------|------------------|------------------|

year 1     year 2

# Thanks!