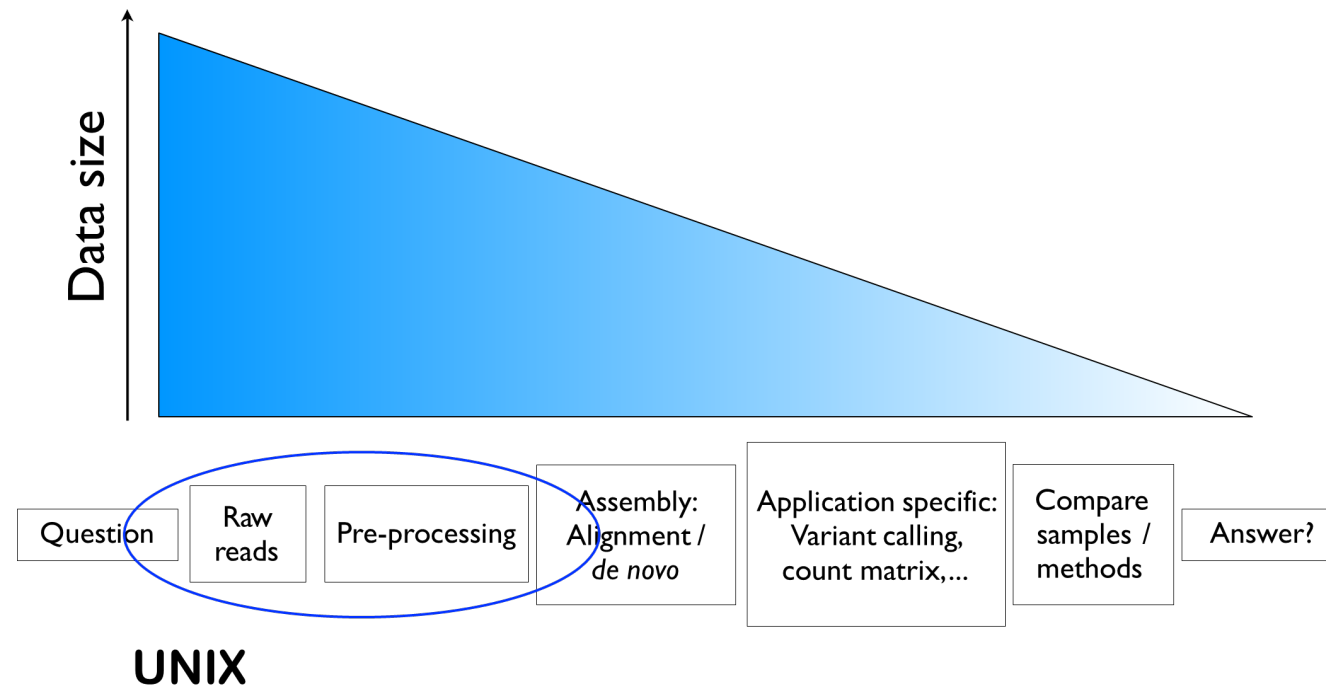**DTU Health Technology**
**Bioinformatics**

# Week 5 Recap

*Gisle Vestergaard*
*Associate Professor*
*Section of Bioinformatics*
*Technical University of Denmark*
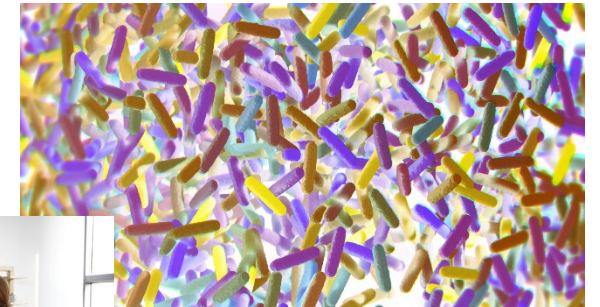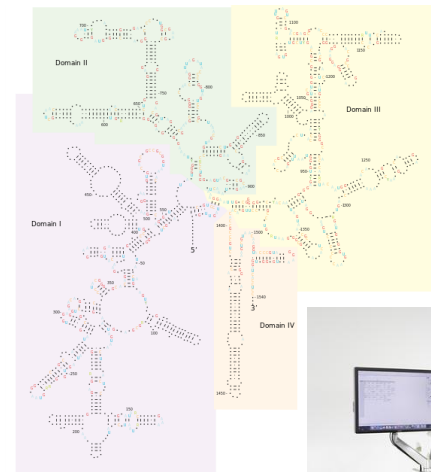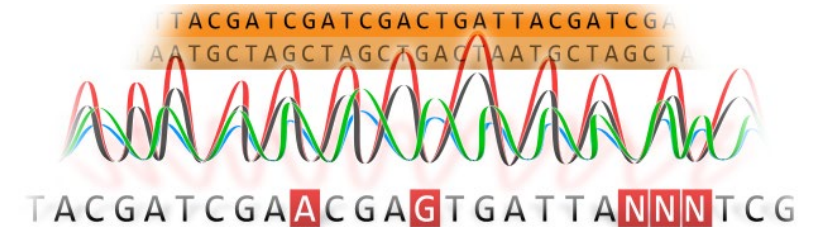*gisves@dtu.dk*

# Last Week

- Getting on the cloud
- Sequencing technology
- Raw sequencing reads
- Preprocessing reads

# Today

- Last Wednesday
- NGS Quiz in Discord groups
- Sequence alignment
- Presentations
- 16s rRNA Amplicon sequencing
- Metagenomics and microbial ecology

Please check webpage for times!

# Technological advances continues to advance metagenomic possibilities

- New Technologies means new possibilities
- …also means new types of errors
- Illumina is the current workhorse
  - Great for many applications
- Long read technology
  - Adding information
  - Resolves difficult regions during genome assembly

# Reminder: things ~~can~~ will go wrong

template
read

mismatch
AGCAATCTCAATTACA**A**ATATACACCAACAAA
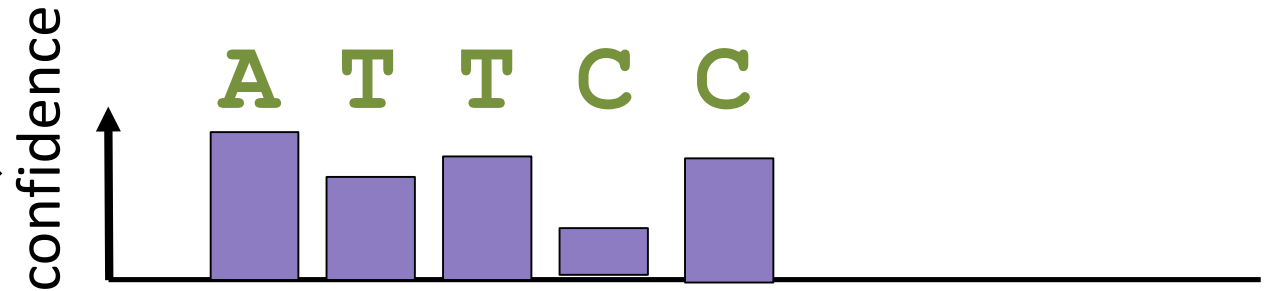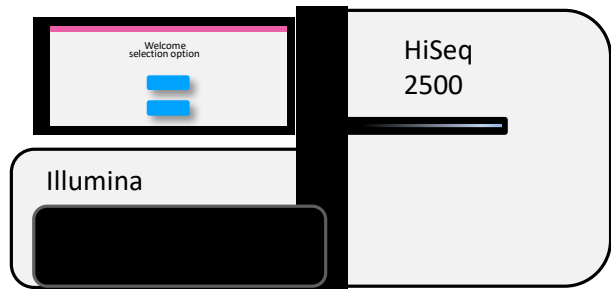AGCAATCTCAATTACA**G**ATATACACCAACAAA

insertion
AGCAATCTCAATTACA**–**AATATACACCAACAA
AGCAATCTCAATTACAC**G**ATATACACCAACAA

deletion
AGCAATCTCAATTACA**A**ATATACACCAACAAA
AGCAATCTCAATTACA**–**ATATACACCAACAAA

# Fastqc reports

- Report basic statistics on your data
- Identify issues with your data
- **Use at each step of preprocessing to check progress**



FastQC Report

## Summary

✅ Basic Statistics
✅ Per base sequence quality
✅ Per sequence quality scores
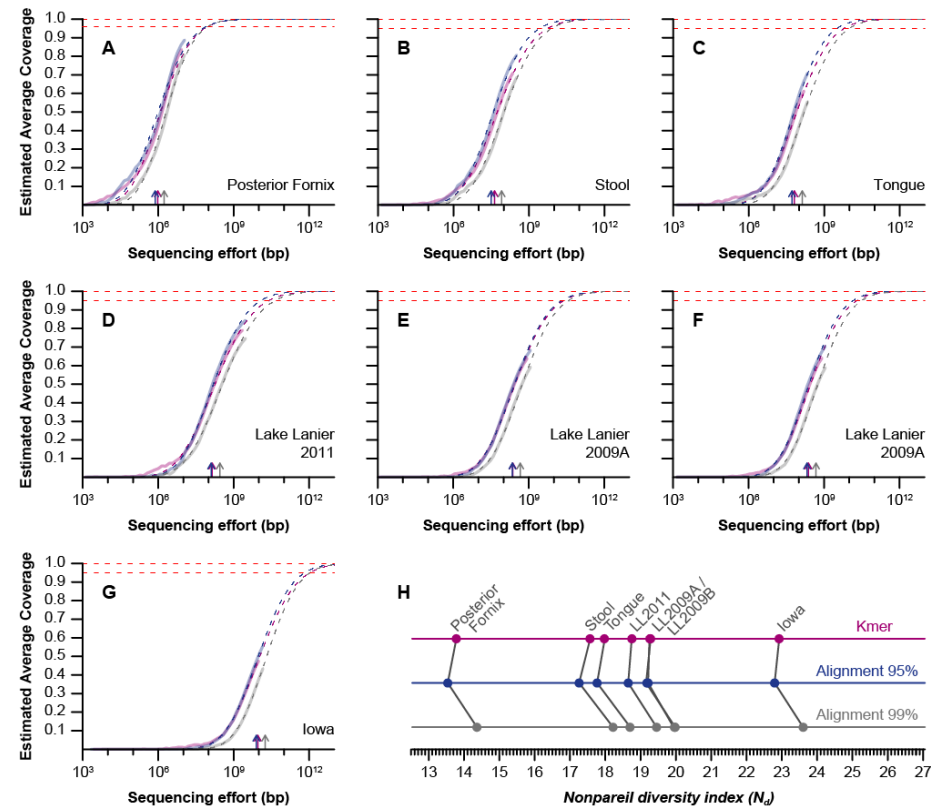✅ Per base sequence content
✅ Per base GC content
⚠️ Per sequence GC content
✅ Per base N content
✅ Sequence Length Distribution
✅ Sequence Duplication Levels
✅ Overrepresented sequences
⚠️ Kmer Content

8

# Sequencing depth for shotgun metagenomics

- No reference database like 16s, therefore we cannot use rarefaction
- Nonpareil: How often do I find the same read in a dataset?

# QUIZ TIME!