

Swiss Institute of Bioinformatics

INTRODUCTION TO SEQUENCING DATA ANALYSIS

### Quality control + database retrieval

Deepak Tanwar Frédéric Burdet

April 23-25, 2025 Adapted from previous year courses



### Why Quality control?

- 1. How is the base quality?
- 2. What is the read length?
- 3. Are there adapters/barcodes in my sequences?
- 4. Are there overrepresented sequences?



### **Dedicated software**

- Manufacturer Tools
  - Built-in QC during base calling and run monitoring
- Commonly Used QC Software:
  - Illumina:
  - FastQC General QC (per-base quality, GC content, adapter detection)
  - MultiQC Aggregates FastQC and other reports for batch analysis
- Oxford Nanopore (ONT):
  - pycoQC Run-level metrics from basecalling (yield, read length, quality)
  - NanoPlot Read length and quality visualization (FASTQ, BAM, etc.)
- ONT & PacBio:
  - NanoStat Summary statistics from long-read sequencing files



### fastq

4

fasta + basequality (fasta + q = fastq)

$$BASEQ = -10log_{10} \Pr\{base \ is \ wrong\}$$

$$-10log_{10} (0.01) = 20$$
  
$$-10log_{10} (0.1) = 10$$
  
$$-10log_{10} (0.5) = 3$$



Phred score



Quality scores across all bases (Sanger / Illumina 1.9 encoding)





% Adapter





## Trimming

Find and remove:

- >> Regions or reads with low base quality
- >> Adapter sequences
- >> poly G sequences (e.g. with NovaSeq 6000)

Software: fastp (or cutadapt, trimmomatic, trim\_galore, bbduk ..)

Articles on frequently occurring issues: <u>sequencing.qcfail.com</u>





# Quality trimming with fastp

Default:

- Remove reads with >40% bases <Q15</li>
- Trim poly N (and poly G)
- Autodetect adapters in R1, for both:
  - --detect\_adapter\_for\_pe





### Databases





### What is INSDC?

- INSDC = International Nucleotide Sequence Database Collaboration
- Founded to ensure free and open access to nucleotide sequence data worldwide
- Three partners:
  - DDBJ (Japan, Asia)
  - ENA (Europe)
  - GenBank (USA, North America)
- Synchronized daily to maintain a shared global repository





Publications PubMed

OTHER DATASETS BioSample

GEO DataSets

GEO Data Details

Data volume, Supplementary Mbytes

Parameter

SRA Data Details

Data volume, Tbytes

Parameter Data volume, Gbases

An offic	ial website of the United States government Here's	<u>s how you know</u> ∽					
NIH	National Library of Me National Center for Biotechnology Info	edicine ormation				Log in	
BioProject  Advanced Browse by Project attributes							
Display Setting Multi-omic heterogene	<sup>gs:</sup>	tic islet donors reveals Accessio pe 2 diabetes (human)	n: PRJNA690574	Send to: - Related information - SJNA690574 ID: 690574 BioSample - Genome		۲	
To gain insi we collected More	ght into the history of islet cell deterioration al d surgical pancreatic tissue samples from 133	ong the progression from normal glycemic regulation to T metabolically phenotyped pancreatectomized patients (PI	C2D, See ( PP). Inform Homo	Genome nation for sapiens	GEO DataSets PubMed SRA Taxonomy		
Accession Data Type Scope	PRJNA690574; GEO: GSE164416 Transcriptome or Gene expression Multiisolate		Navigat 88717 projects by or	re Across additional are related ganism.			
Organism	Homo sapiens [Taxonomy ID: 9606] Eukaryota; Metazoa; Chordata; Craniata; Vertebra Haplorrhini; Catarrhini; Hominidae; Homo; Homo s	ata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primate sapiens	25;		Multi-omics profiling of livin pancreatic islet donors revo	Turn Off Clear ig human eals heten BioProject	
Publications Submission	Wigger L et al., "Multi-omics profiling of living cell trajectories towards type 2 diabetes.", Na Registration date: 7-Jan-2021	human pancreatic islet donors reveals heterogeneous bet it Metab, 2021 Jul;3(7):1017-1031	a		Q SRP021519 (8)	SRA	
Relevance	Medical				<ul> <li>Entrez Direct: E-utilities on Command Line - Entrez Pr</li> <li>SRR519926 (1)</li> </ul>	the Unix ogramming	
Resource Name         Number of Links           SEQUENCE DATA         Very Sequence Data		Number of Links				See more	
SRA Experiments		133					

1

133

Value

36

Value

349

0.23

1



#### & DDBJ · Services SuperComputer Statistics Activities About Us

#### home > bioproject > PRJNA690574

identifier	PRJNA690574				
type	bioproject				
sameAs	GEO	<u>GSE164416</u>			
organism	Homo sapiens				
title	Multi-omics profiling of living human pancreatic islet donors reveals heterogeneous beta cell trajectories toward type 2 diabetes				
description	To gain insight into the history of islet cell deterioration along the progression from normal glycemic regulation to T2D, we collected surgical pancreatic tissue samples from 133 metabolically phenotyped pancreatectomized patients (PPP). Gene expression profiles of islets isolated by laser capture microdissection (LCM) from resected and snap-frozen pancreas samples were assessed by RNA sequencing.Overall design: This study includes RNA-Seq samples from pancreatic c islets of 133 human donors, stratified into four groups based on their diabetes status: 18 were non-diabetic (ND), 41 had impaired glucose tolerance (IGT), 35 had Type 3c diabetes (T3cD), and 39 had Type 2 diabetes (T2D). The group a ssignments are based on thresholds defined in the guidelines of the American Diabetes Association.For data analysis, a subset of 92 pancreatic islet samples was defined, which included only those samples in which the gene INS showed the highest expression (i.e., highest normalized counts value). Statistical analyses were performed both on the complete transcriptomics data set and on this restricted data set.				
data type	Transcriptome or Gene expression				
organization					
publication	34183850				
external link					



### For sensitive human data

# EUROPEAN GENOME-PHENOME ARCHIVE

See also: https://ega-archive.org/about/projects-and-funders/federated-ega/



15

#### 16

### EGA

- EGA = European Genome-phenome Archive (EMBL-EBI, Europe)
- Designed for controlled access to human data with privacy concerns
- Ideal for:
  - Clinical studies
  - Patient phenotypes
  - Genomic variants
- Access requires data access committee (DAC) approval



### Summary

17

- INSDC enables global sharing of nucleotide sequence data
- GEO / SRA are part of this ecosystem, focusing on expression/functional data
- For sensitive human datasets, EGA is the recommended platform
- Choosing the right repository ensures data availability and compliance



### Command line tools

Retrieve raw data: SRA-tools >> prefetch >> fastq-dump

Retrieve sequences: Entrez Direct >> esearch >> efetch





### Thank you

DATA SCIENTISTS FOR LIFE





sib.swiss