



Quantitative Molecular Neurogenomics 2023

Module 5 - Data Sources, Acquisition, & QC

ME.440.825

Loyal Goff - Course Director

Genevieve Stein-O'Brien - Instructor

Richard Sriworarat - TA

Kyla Woyshner - TA

Jeanette Johnson - TA

October 2nd, 2023

PUBLIC NGS DATA SOURCES

NGS Data Sharing is required for NIH-funded research

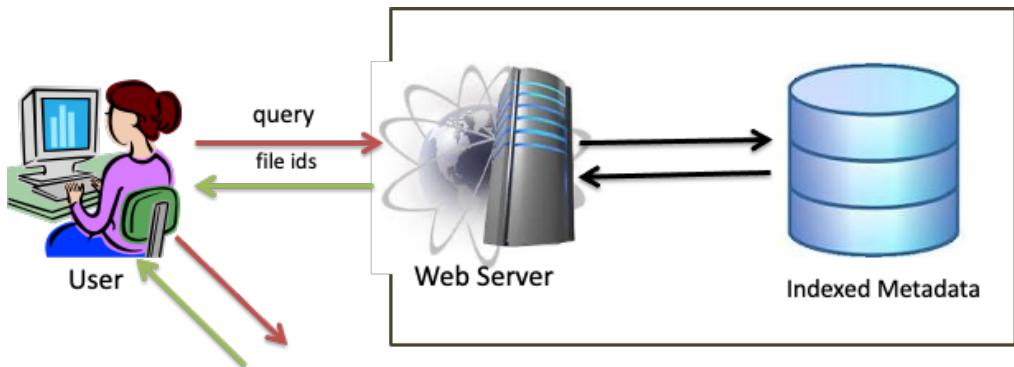
- NIH Issued Genomic Data Sharing (GDS) policy, August 2014
 - “...applies to all NIH-funded research (e.g., grants, contracts, and intramural research) that generates large-scale human or non-human genomic data.”
 - “NIH strongly encourages the use of established repositories to the extent possible for preserving and sharing scientific data.”
 - “Shared scientific data should be made accessible as soon as possible, and no later than the time of an associated publication”

Common NGS Data sources/archives

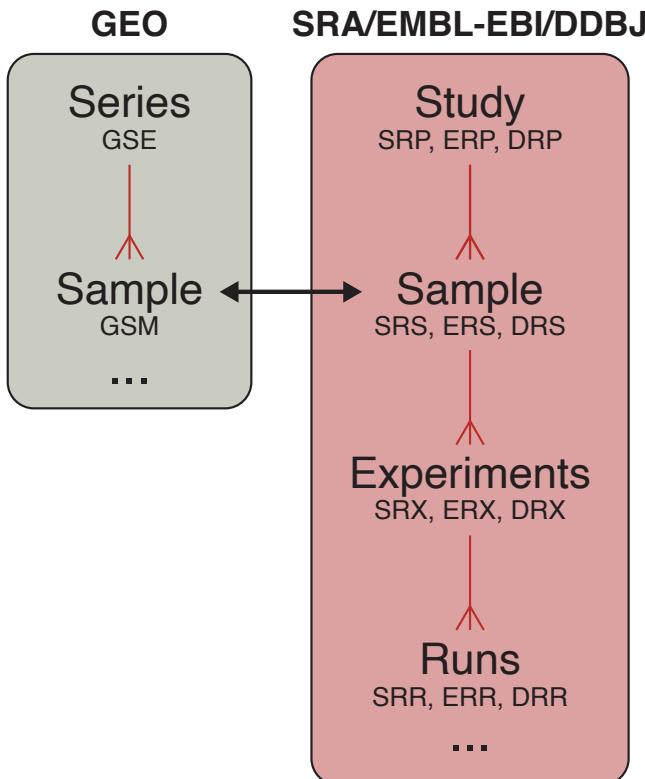
- Gene Expression Omnibus (GEO)
 - Public functional genomics data repository
 - Storage point for ***processed*** data files
- Sequence Read Archive (SRA)
 - Stores ***raw*** sequencing reads and associated metadata
- EMBL-EBI
 - European Bioinformatics Institute
- DDBJ
 - DNA Data Bank of Japan

SRA - ncbi.nlm.nih.gov/sra/

- Archives raw sequencing data and alignment information from high-throughput sequencing platforms
- Data are ingested from direct submissions to GEO
 - Processed into ‘SRA’ format for efficient storage



Navigating Common NGS Primary Data Archives



- Series/Studies:
 - Single Accession for all downstream elements of a project
- Samples:
 - A distinct biological sample
 - Usually 1:1 across Data sources
- Experiment:
 - A unique sequencing library and sequencing method for a specific sample (sample + assay = library/experiment)
- Run:
 - The execution of an experiment on a sequencing platform and associated data (output)
 - Different runs of the same experiment can *usually* be aggregated together
 - (e.g. same library sequenced twice or on separate flow cells)

What am I looking for...

- Raw files for a given experiment can/should include:
 - Raw sequencing output (.fastq)
 - Sample information/metadata
 - Including how the raw fastq files relate to samples/replicates
 - Processed data files
 - .csv/.tsv/.txt/.mtx
 - These files store the ‘processed values’ that the authors used for downstream analyses.
 - Any specialized accessory data
 - e.g. Custom references or annotations

DATA ACQUISITION



≡ HOME MAGAZINE COMMUNITY INNOVATION

NEWSLETTER ABOUT

SUBMIT MY RESEARCH

LOG IN/REGISTER



Tools and Resources

Neuroscience

Hipposeq: a comprehensive RNA-seq database of gene expression in hippocampal principal neurons

Download

Cite

Comment



19,126 views

189 citations

Mark S Cembrowski, Lihua Wang, Ken Sugino, Brenda C Shields, Nelson Spruston

Janelia Research Campus, Howard Hughes Medical Institute, United States

Apr 26, 2016 · <https://doi.org/10.7554/eLife.14997>

<https://elifesciences.org/articles/14997>

ACCESSORY DATA SOURCES

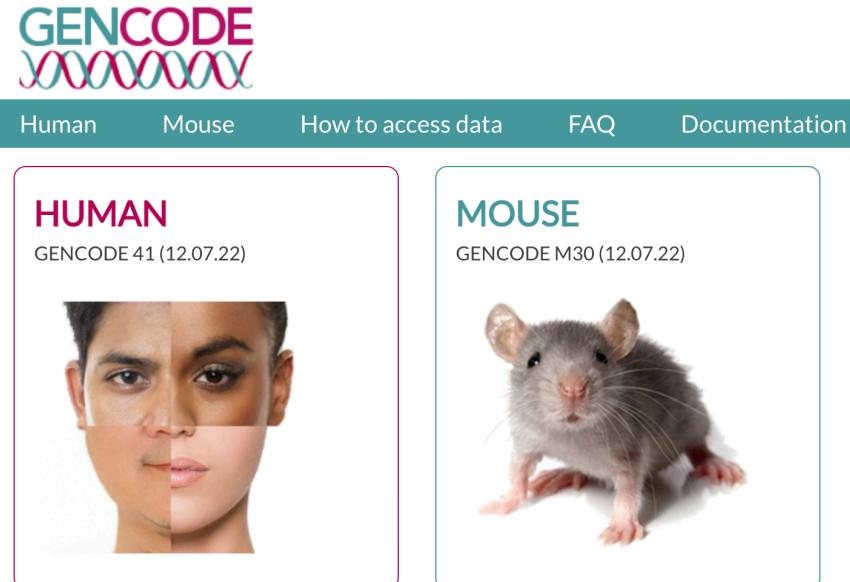
Reference Genomes

- Most applications will require a reference genome sequence against which to map
- Genome assembly/annotation projects are continuously being updated
 - Released in tagged versions
- Sources:
 - www.ensembl.org
 - genome.ucsc.edu
 - Consortia/species-specific sites
 - flybase.org
 - rgd.mcw.edu (Rat Genome Database)

The screenshot shows the Ensembl genome browser homepage. At the top, there's a dark header bar with the Ensembl logo and links for BLAST/BLAT, VEP, Tools, BioMart, Downloads, Help & Docs, and Blog. Below the header, there are four main navigation sections: 'Tools' (with a 'All tools' link), 'BioMart >' (described as 'Export custom datasets from Ensembl with this data-mining tool'), 'BLAST/BLAT >' (described as 'Search our genomes for your DNA or protein sequence'), and 'Variant Effect Predictor >' (described as 'Analyse your own variants and predict the functional consequences of known and unknown variants'). The main content area has a light blue background. It features a search bar with dropdown menus for 'All species' and 'for', and a 'Go' button. Below the search bar is a note: 'e.g. BRCA2 or rat 5:62797383-63627669 or rs699 or coronary heart disease'. Further down, there's a section for 'All genomes' with a dropdown menu labeled 'Select a species'. To the right of this is a 'Favourite genomes' section with icons for Human (GRCh38.p13), Still Using GRCh37?, Mouse (GRCm39), and Zebrafish (GRCz11). At the bottom left, there's a link 'View full list of all species'.

Reference Transcriptome(s)

- GENCODE
 - Integrated annotation of gene features
 - gencodegenes.org
 - (Current) highest-quality source for transcriptome sequence and annotation in human and mouse
- Ensembl aggregates many other species assembled transcriptomes



Reference Gene/Feature Annotations

- Common formats:
 - GFF/GTF
 - <https://useast.ensembl.org/info/website/upload/gff.html>
 - Holds information about gene structure, ids, position
 - General Feature Format (.gff)
 - General Transfer Format (.gtf)
 - BED (Browser Extensible Data)
 - <http://genome.cse.ucsc.edu/FAQ/FAQformat.html#format1>
 - Generalized file format for genomic regions as coordinates and associated annotations

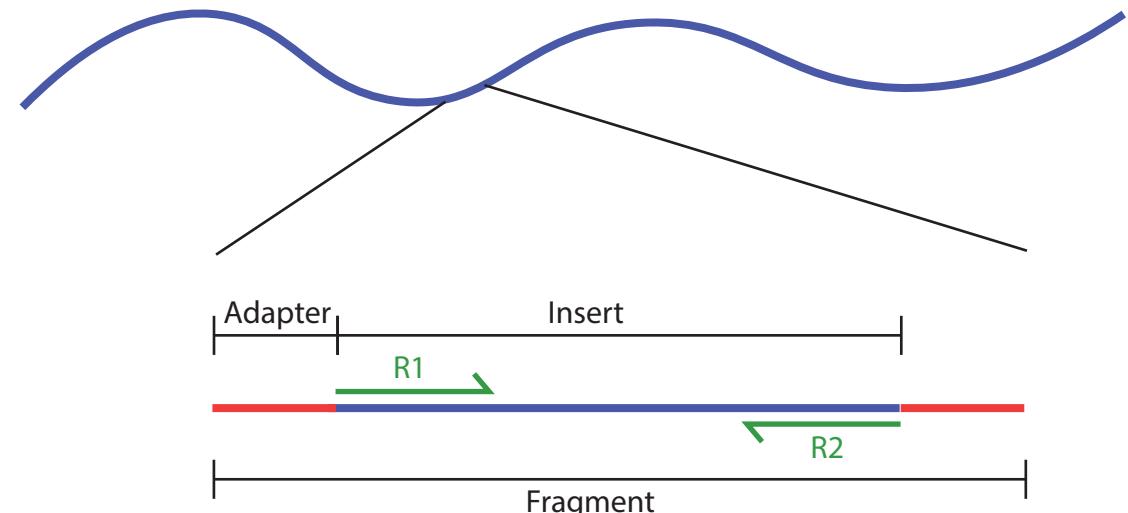
Building your own genome/transcriptome assembly

- Assembly of references is a massive and challenging undertaking.
- Is only ***required*** if:
 - You are working in a non-standard model organism
 - You anticipate novel genes/gene structures/features to be present in your dataset that wouldn't have been previously discovered/annotated
- Often, supplements are needed to assemblies/annotations:
 - Adding a BAC sequence that is part of a transgenic line
 - Adding transcript information for GFP in a reporter assay

WHAT IS A SEQUENCING READ?

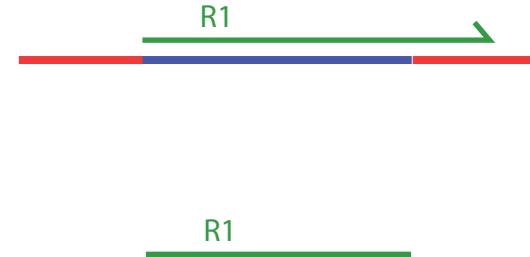
Read terminology

- Input material is *fragmented*
- *Fragments* have *adapter* sequences added to either end
 - PCR handles
 - Sequencing adapters
 - Index/barcode sequences
- *Insert* is the portion of the *fragment* containing the sequence of interest
- *Reads* are the portions of the *insert* that are read by the sequencer



Adapter Trimming

- Why would we need to trim?
 - Shorter fragments and/or longer reads run the risk of actually sequencing *into* the downstream adapter sequences
 - Additionally, bases at the ends of reads tend to have lower quality scores
- This **can** have negative consequences for finding *alignments* to a reference
- Many tools exist to test the ends of reads for adapter sequences or lower quality bases and ‘trim’ them.
- Choice of whether to trim is subjective



Common Read Adapter Trimming Tools

- Trimmomatic (<http://www.usadellab.org/cms/?page=trimmomatic>)
- FASTX-Toolkit
- fastp
- Cutadapt (<https://cutadapt.readthedocs.io/en/stable/>)

Single End vs Paired End sequencing

- Fragments can be sequenced from one OR both ends
 - Decision made at sequencing stage
 - Aligners require different input information for each
- Paired-end is almost **always** preferred
 - Significantly better alignment rates
 - Empirically determine fragment length (a useful metric for quantification)
 - Double the cost for reads of the same length
 - Read pairs are aligned jointly
 - ‘Properly’ mapped pairs significantly reduce ambiguity in alignment

Single End Sequencing



<sample_name>.fastq.gz

Paired End Sequencing



<sample_name>_1.fastq.gz

<sample_name>_2.fastq.gz

READ QC

Why do we need to QC reads?

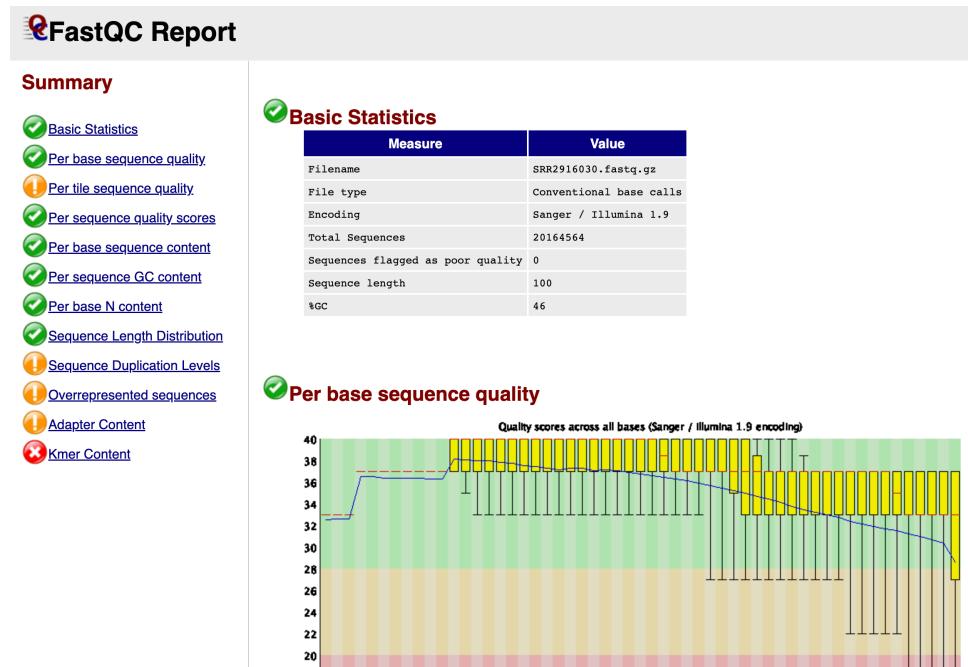
- Ensuring data integrity
 - Different samples/labs/assays have different efficiencies and biases
 - Need reads of comparable quality to make fair comparisons across replicates and conditions.
- Poor quality reads can affect mapping rates, and/or lead to misquantifications or misinterpretations
- Identifying and filtering out low quality reads can improve analysis even if you have to sacrifice some of the data.

Read QC - What are we looking for?

- Sufficient ***number*** of reads
 - From sequencing run
 - Aligning to reference
- Reads of ***consistent*** and ***adequate*** quality scores
- Consistent sequence ***content***
 - Per base distribution of nucleotides
 - Per read GC content
- Low number of poor quality or ambiguous bases
- Sufficient library ***complexity***
 - Low number of over-represented sequences
 - Low number of over-represented kmers
- Low proportion of sequencing adapters

fastqc (<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>)

- A GUI & command-line tool to analyze .fastq files for common quality metrics
- fastqc <fastq_file.fastq.gz>
- fastqc *.fastq.gz
- echo \$(ls *.fastq.gz) | xargs -n 1 -P 4 fastqc



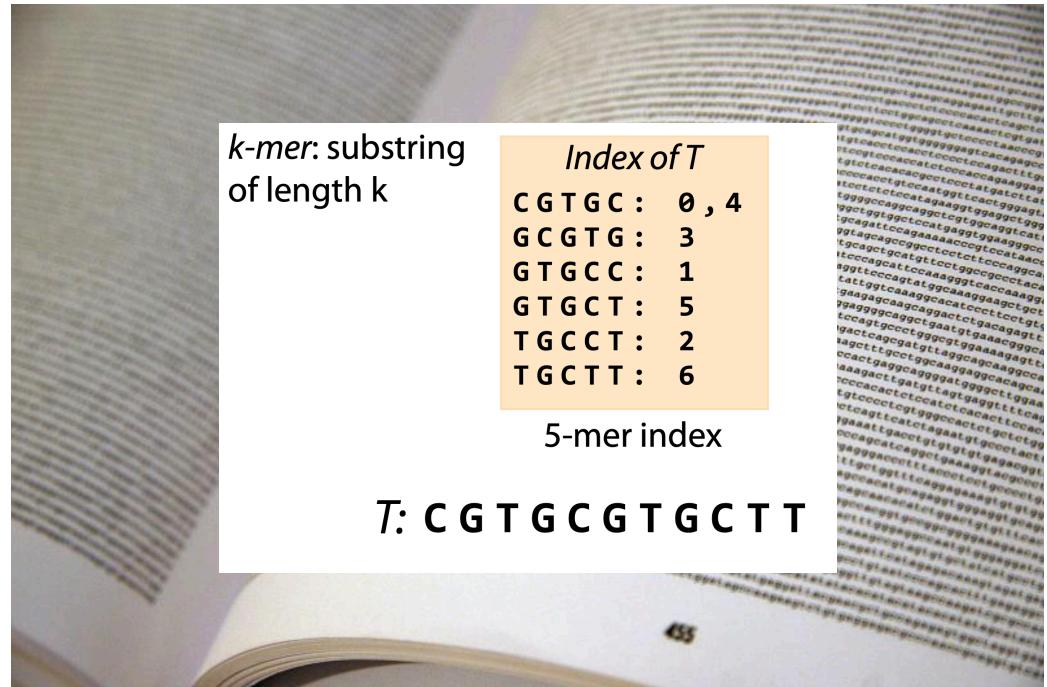
MAPPING/ALIGNMENT

Read alignment

- Goal: to locate the genome coordinates from which a sequenced read could originate
- Requires searching a reference (genome/transcriptome) to find sequence matches for each of millions of reads
 - With some mismatch tolerance (SNPs, Indels, sequencing errors)

Reference Indexing

- Searching a linear reference sequence for substring matches is a computationally intensive task
 - Printed human genome, single spaced, 12pt font, 1" margins requires 1,206,980 pages (http://bio4.us/biotrends/human_genome_height.html)
- Generating an index (usually constructed from Kmer ‘seeds’) helps aligners to jump to specific regions of the genome that can then be extended to a full read match
- Most modern aligners provide a mechanism/software to generate an indexed reference as part of their documented workflows.



<https://www.langmead-lab.org/teaching.html>

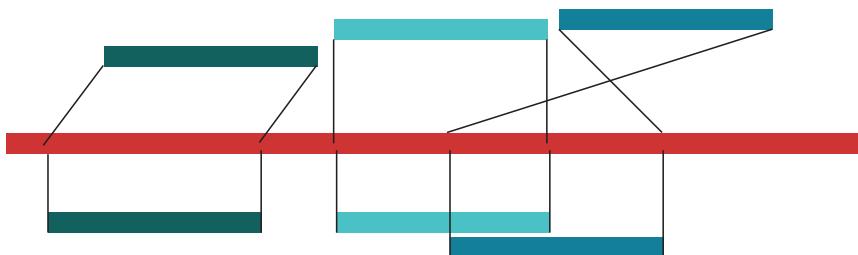
https://github.com/BenLangmead/ads1-slides/blob/master/0300_index_index.pdf

Choosing the appropriate reference

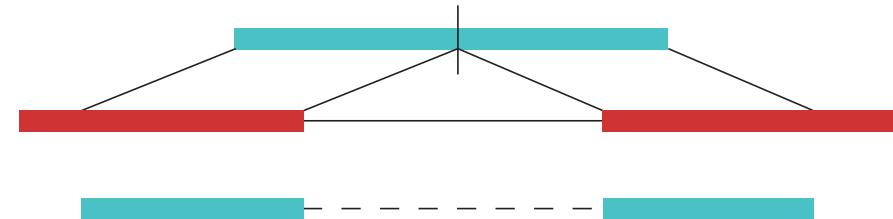
- Primarily guided by the 'purpose' of your analysis
 - Differential expression
 - Variant calling
 - Structural variant
 - Peak discovery
- Version and annotation level
 - Aim to use the most recent version that has sufficient annotation for your downstream needs
- Specialized/modified references
 - Ex: transcriptomes for RNA-Seq, mitochondrial genomes, or plasmids for specific studies.
 - Adding transgenes or other exogenous sequences included in your study to your reference may be helpful.
- Closeness to your sample
 - Ideally, use a reference from the same species/subspecies/strain.
 - Consider using a closely related species' genome if studying a less-known organism.

Ungapped vs Gapped alignment

Ungapped/unspliced



Gapped/spliced

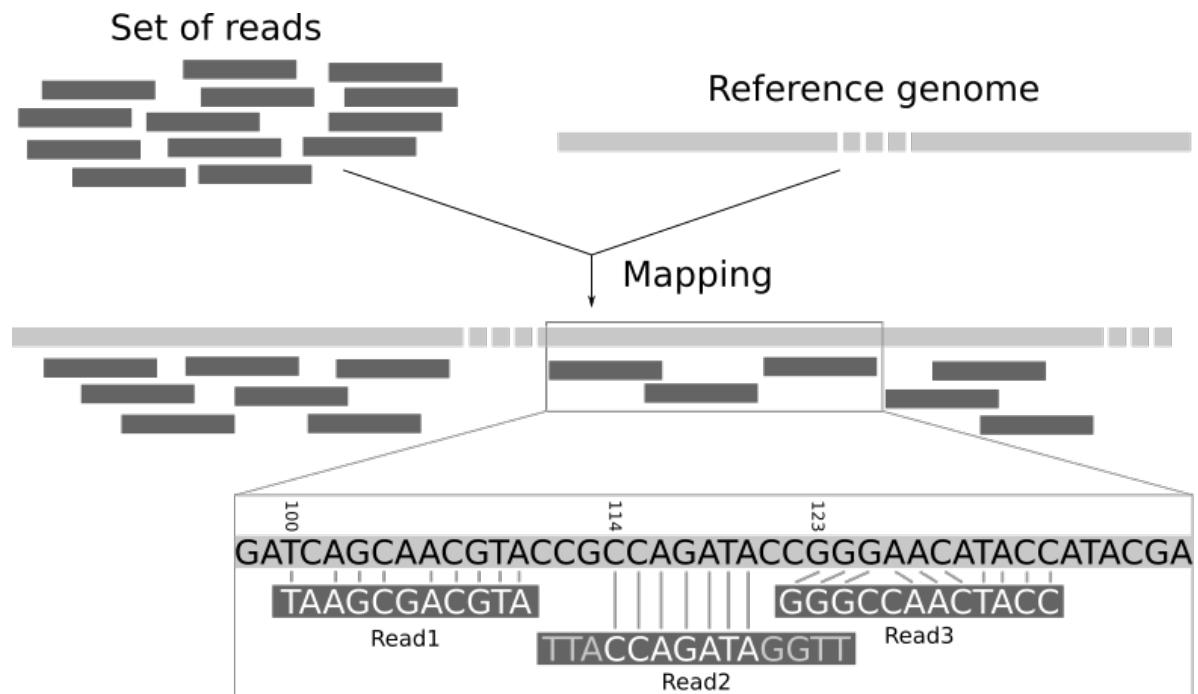


- Alignments without insertions, deletions, or gaps.
- **Appropriate** for most genomic reads
 - e.g. ChIP-Seq, ATAC-seq
- Speed: Faster, as there are fewer alignment possibilities.
- Ideal for short sequences or when expecting high similarity.
- Limitation: Cannot capture some structural events like insertions or deletions.

- Alignments that allow for gaps, representing possible insertions, deletions in read relative to reference
- **Required** for most RNA-Seq studies
 - Splicing produces exon-exon junctions that can only be mapped with gaps to contiguous reference
- Accuracy: More accurate representation of evolutionary events, indels, mutations, splicing, etc
- Essential for longer sequences or when inferring evolutionary relationships.
- Limitation: Computationally more intensive.

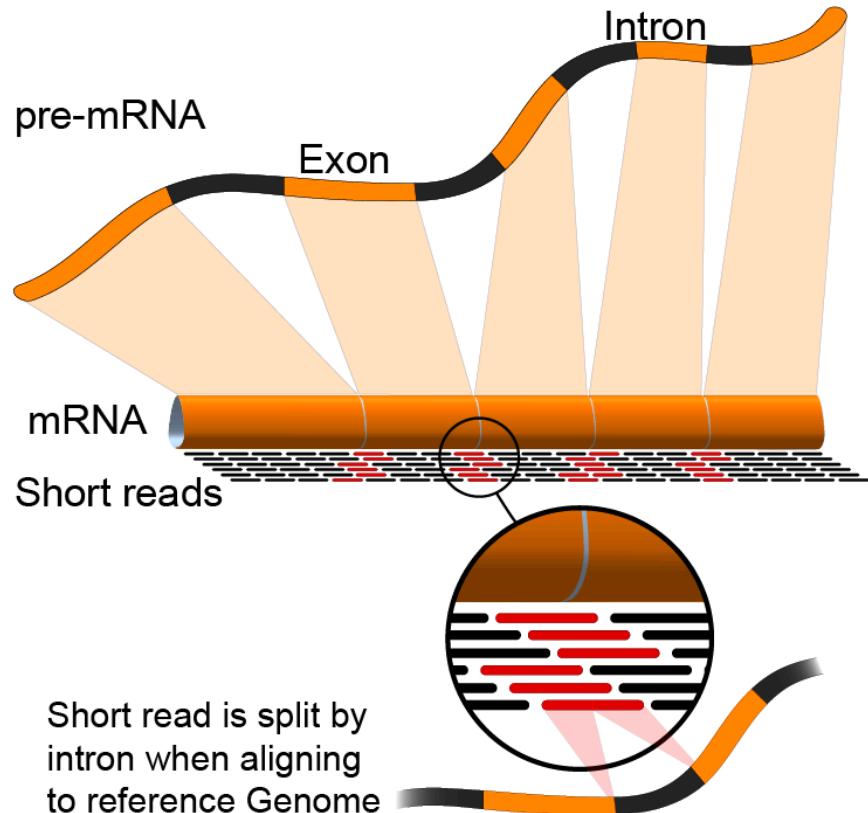
Alignment of DNA fragments to genome (ATAC-Seq, ChIP-Seq, etc)

- Each DNA-based read is assigned to a specific location in a reference genome
- Peaks identified by boundaries of overlapping reads above features of interest
 - Peak height/width/ count all useful attributes for quantification
 -



Alignment of RNA-based fragments

- Reads derived from cDNA may contain exon-exon junctions
 - Longer & paired-end reads more useful in this context
- Require use of splice-aware aligner
- Junctions useful for identification & quantification of isoform-level expression



What aligner to use?

- Depends on:
 - Task
 - Data type
 - Biological questions
 - Available resources
- *Tip: For choosing alignment tool and parameters for your final projects, check methods section, supplement, and data archive records for description of prior efforts.*

Common Aligners

- Unspliced (ungapped) Aligners
 - MAQ
 - BWA
 - Bowtie2
- Spliced (gapped) Aligners
 - Tophat2
- Pseudoalignment
 - Kallisto
 - Salmon
- HiSat2
- STAR

How can we improve the speed/efficiency of this process?

- More efficient algorithms:
 - Indexed search of genome/transcriptome
- Parallelization
 - Multiple threads for each iteration/call of the aligner
 - Multiple concurrent job submissions across samples
 - Naive parallelization is one chief benefit of HPC systems
 - Jobs are submitted to a ‘scheduler’ which balances resources and submits
 - Slurm (<https://slurm.schedmd.com/documentation.html>) is available on rockfish

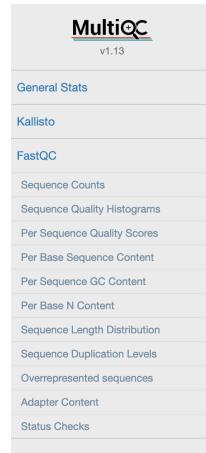
ALIGNMENT QC

Alignment QC - What are we looking for?

- **Number** of reads aligned per sample
- **Percentage** of reads aligned per sample
- **Number/percentage** of reads uniquely mapping
 - % of reads with multiple mapping to reference (ambiguous)
- Paired End read **concordance**
 - # of PE reads uniquely mapping
 - # of PE reads mapping but not together (discordant)
 - # of PE reads where only one mapped [uniquely]
- **Consistent** scores/values across samples
- *Values are often reported after each sample is processed. Useful to try and capture this information in a log file.*

Aggregate QC Reports with MultiQC

- <https://multiqc.info/>
- Generates a single .html summary report across multiple stages of a workflow
- Can be called (multiple times) from the root directory of a project
 - multiqc .
- Parses outputs and logs of many common genomics tools
 - Often times, STDERR output contains needed information and should be redirected to a file using (2> <my_job>.log)



A modular tool to aggregate results from bioinformatics analyses across many samples into a single report.

Report generated on 2022-09-22, 20:49 EDT based on data in: /Volumes/GoogleDrive/My Drive/Work/Goff Lab/Teaching/Quantitative Neurogenomics/01



General Statistics

Sample Name	% Aligned	M Aligned	% Dups	% GC
SRR2916027	42.7%	14.9	33.4%	43%
SRR2916028	44.3%	17.7	33.0%	44%
SRR2916029	46.2%	16.0	36.6%	44%
SRR2916030	47.0%	9.5	38.1%	46%
SRR2916031	43.8%	8.4	34.5%	46%

ALIGNMENT VISUALIZATION WITH IGV

Integrative Genomics Viewer (IGV)

- <https://software.broadinstitute.org/software/igv/>
- Useful to generally inspect alignments where possible or relevant
 - Genomic deletion sites
 - Mutation insertion sites
 - Key regions/genes relevant to study
- Import indexed .bam files (and references) to view read distributions and associated features
 - Identified exon-exon junctions

