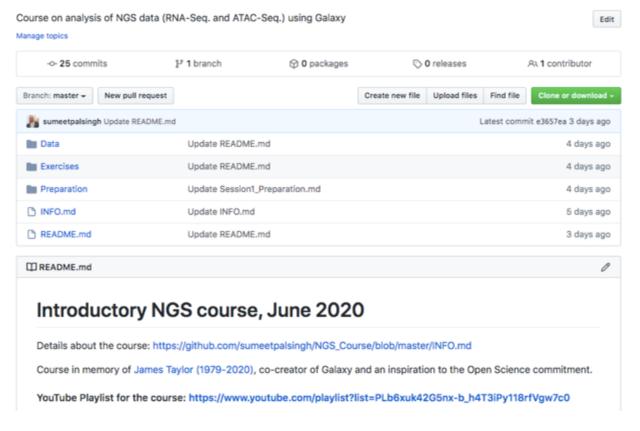# NGS Course

Session 1

Sumeet Pal Singh and Yura Song

# NGS Course

- Session 1 (05 June)
  File types in NGS (fastq, sam / bam, genome index and gtf / gff3)
  Mapping Fastq files

- Session 2 (18 June)
  Designing and saving workflow in Galaxy
  RNA-Seq. analysis
  Controlling for covariates in RNA-Seq. analysis

- Session 3 (25 June)
  ATAC-Seq. pipeline
  Interrogating ATAC-Seq. data for peaks, tf binding sites, enriched motifs

# Course Repository

- https://github.com/sumeetpalsingh/NGS_Course

# Post Issues on Github

# What to expect from the course

- Using Galaxy

- Allow you to go from Raw Data to Analysis

- Teach the steps and associated tools for analysis pipeline

- Independent to work with NGS-Data

- Develop Galaxy Workflows for new analysis / pipeline

# What the course is not about

- Does not cover tools / pipelines not present in Galaxy
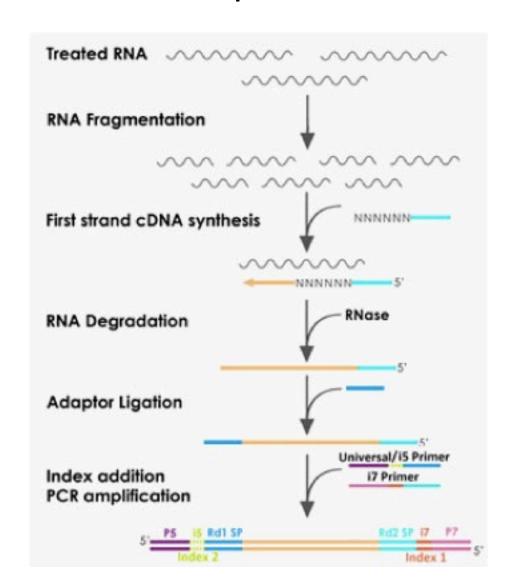- Does not cover executing the tools in Shell Script / HPC
- Does not cover adding features to tools that are not implemented in Galaxy
- Does not cover working with non model-organisms

- Only covers bulk RNA-Seq. and bulk ATAC-Seq. analysis (not single-cell) made using Illumina instrument

# RNA-Seq.

# RNA-Seq. Library Preparation



https://www.youtube.com/watch?v=-kTcFZxP6kM

# RNA-Seq. Library

# Flow cell-based Sequencing

# Flow cell capacity

## Reads Passing Filter Per Flow Cell

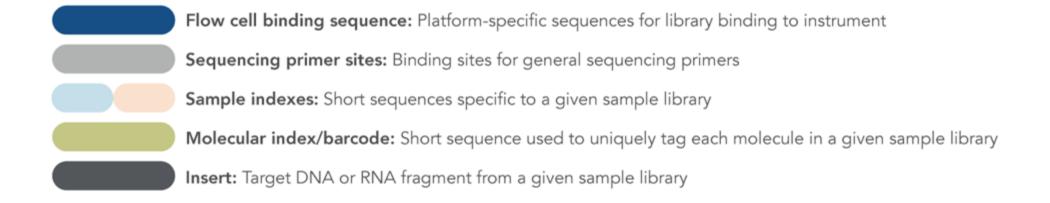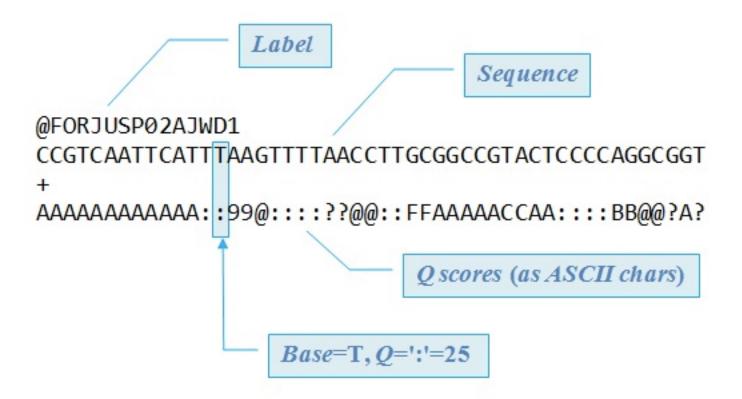| | NovaSeq 6000 System | | | |
|---|---|---|---|---|
| Flow Cell Type | SP | S1 | S2 | S4 |
| Single-end Reads | 650–800 M | 1.3–1.6 B | 3.3 B–4.1 B | 8-10 B |
| Paired-end Reads | 1.3–1.6 B | 2.6–3.2 B | 6.6–8.2 B | 16–20 B |

For a regular RNA / ATAC-Seq. sample: 10 – 50 M

# Raw Data

- Fastq format: Fasta format with quality scores

Fasta



| Header | >VIT_201s0011g03530.1 |
| Sequence | AATTAAGCATAAATACTCACTCTTACCCCCTTATTTTCTTATCTCTCATCACTTTTGGTGCGAAG |
| | GACCATGAGAACAAGCTGCAATGGGTGTAGGGTTCTTCGCAAGGCATGCAGCCAAGACTGCATCA |
| Header | >VIT_201s0011g03540.1 |
| Sequence | CAGGTAGCGTGAAGTTAAACCCTAGCGCTTTAGACAAACAGCTGTAGTCACCGCCCACAAACACC |
| | AGCCTCTGAGACACCACCTCAAACCTTTCCACTTAAATACACATCCCTCACACCCTTTTCAATTC |
| Header | >VIT_201s0011g03550.1 |
| Sequence | CATGCAAAGCTGAACGCGATGCTGTGATTGGTGGTAAGTGGTAGTTGAGTAAATTTGACAGTGAA |
| | GCCGAAATGGTAAAAGACTAAGGCTAGAAGTAGAATACCACTGTTCTTCTCATCACGTGGGCCCA |

# Raw Data

- Fastq format: Fasta format with quality scores



Label

Sequence

@FORJUSP02AJWD1
CCGTCAATTCATTTAAGTTTTAACCTTGCGGCCGTACTCCCCAGGCGGT
+
AAAAAAAAAAAA::99@::::??@@::FFAAAAACCAA::::BB@@?A?

Q scores (as ASCII chars)

Base=T, Q=':'=25

# Quality Score

| | |
|---|---|
| Clipped length: | 663 |
| Left clip: | 14 |
| Right clip: | 676 |
| Avg. qual. in clip.: | 46.15 |

| | |
|---|---|
| Samples: | 16978 |
| Bases: | 733 |
| Average spacing: | 24.0 |
| Average quality >= | 10: 28, 20: 108, 30: 548 |

| Quality: | |
|---|---|
| 0 - 9 | |
| 10 - 19 | |
| 20 - 29 | |
| >= 30 | |

Page: 1 / 3
03.06.2020

N N N N N T T N N GN N N NC N C G A A G G A G G C C G C A A C T T G T T T A T T G C A G C T T A T A A T G G T T A C A A A T A A A G C A A T A G C A T C A C A A T T T C A C

A A A T A A A G C A T T T T T T T C A C T G C A T T C T A G T T G T G G T T T G T C C A A A C T C A T C A A T G T A T C T T A T C A T G T C T G G A T C T T A A T T A

A T C A G A C C T C A A G G G A A C C C A G T G T G A T T C C G G C A G C G G T C A C G A A C T C C A G C A G G A C C A T G T G A T C G C G C T T C T C G T T G G G G

# Fastq Quality Scores

| Q | P_error | ASCII | | Q | P_error | ASCII | | Q | P_error | ASCII | | Q | P_error | ASCII | |
|---|---------|-------|---|----|---------|-------|---|----|---------|-------|---|----|---------|-------|---|
| 0 | 1.00000 | 33 | ! | 11 | 0.07943 | 44 | , | 22 | 0.00631 | 55 | 7 | 33 | 0.00050 | 66 | B |
| 1 | 0.79433 | 34 | " | 12 | 0.06310 | 45 | - | 23 | 0.00501 | 56 | 8 | 34 | 0.00040 | 67 | C |
| 2 | 0.63096 | 35 | # | 13 | 0.05012 | 46 | . | 24 | 0.00398 | 57 | 9 | 35 | 0.00032 | 68 | D |
| 3 | 0.50119 | 36 | $ | 14 | 0.03981 | 47 | / | 25 | 0.00316 | 58 | : | 36 | 0.00025 | 69 | E |
| 4 | 0.39811 | 37 | % | 15 | 0.03162 | 48 | 0 | 26 | 0.00251 | 59 | ; | 37 | 0.00020 | 70 | F |
| 5 | 0.31623 | 38 | & | 16 | 0.02512 | 49 | 1 | 27 | 0.00200 | 60 | < | 38 | 0.00016 | 71 | G |
| 6 | 0.25119 | 39 | ' | 17 | 0.01995 | 50 | 2 | 28 | 0.00158 | 61 | = | 39 | 0.00013 | 72 | H |
| 7 | 0.19953 | 40 | ( | 18 | 0.01585 | 51 | 3 | 29 | 0.00126 | 62 | > | 40 | 0.00010 | 73 | I |
| 8 | 0.15849 | 41 | ) | 19 | 0.01259 | 52 | 4 | 30 | 0.00100 | 63 | ? | 41 | 0.00008 | 74 | J |
| 9 | 0.12589 | 42 | * | 20 | 0.01000 | 53 | 5 | 31 | 0.00079 | 64 | @ | 42 | 0.00006 | 75 | K |
| 10 | 0.10000 | 43 | + | 21 | 0.00794 | 54 | 6 | 32 | 0.00063 | 65 | A | | | | |

# Using FTP Client to Transfer Files to FTP Server

# Connect to FTP Server

- For galaxy.org: https://galaxyproject.org/ftp-upload/
  FTP Server: usegalaxy.org


- For galaxy.eu: https://galaxyproject.eu/ftp/
  FTP Server: galaxy.uni-freiburg.de


- For galaxy.au: https://usegalaxy-au.github.io/posts/2019/03/18/new-ftp-upload-url/
  FTP Server: usegalaxy.org.au

# Genome Index

- Genome Index <-> Genome
- Dictionary <-> Words

# Genome Index



Suffix tree

Suffix array

Genome Index is unique to every mapping tool

# Mapping Tools

- Bowtie2: http://bowtie-bio.sourceforge.net/bowtie2/index.shtml

- BWA: http://bio-bwa.sourceforge.net/

- Hisat2: https://ccb.jhu.edu/software/hisat2/manual.shtml (Normal Laptop)

- STAR: https://github.com/alexdobin/STAR (High RAM requirements: Human ~32 Gb)

# SAM file format (Alignment Formats)

- SAM – Sequence Alignment / Map Format

- Plain Text (Human Readable)

- Contains
  Quality Scores, Sequence info (Fastq) +
  Alignment Info + MetaData

**HEADER** containing metadata (sequence dictionary, read group definitions etc)
**RECORDS** containing structured read information (1 line per read record)

# SAM Format Example

Chromosome (Mapped database) information

Used program and its variables

Mapped read in forward direction on Chr5

```
@SQ      SN:Chr1  LN:30427671
@SQ      SN:Chr2  LN:19698289
@SQ      SN:Chr3  LN:23459830
@SQ      SN:Chr4  LN:18585056
@SQ      SN:Chr5  LN:26975502
@PG      ID:bwa   PN:bwa    VN:0.5.9-r16
SRR038985.100      0          Chr5      22828962            37          33M        *
0        0          GCCGGTGATGTAATCAAAATATTTGCTACTCTT            WZYTWWTW\]
YVUOW]OEKNUUX]PJSRY][63           XT:A:U   CM:i:0   X0:i:1   X1:i:0   XM:i:
1   XO:i:0   XG:i:0   MD:Z:33
SRR038985.200      0          Chr3      14197678            0           33M        *
0        0          ACCTGGTTGATCCTGCCAGTAGTCATATGCTTG            X]]KN]]
YWUX]XIKYRCHSUYX[[SNQJL[MO        XT:A:R   CM:i:0   X0:i:2   X1:i:0
XM:i:0   XO:i:0   XG:i:0   MD:Z:33 XA:Z:Chr2,+3707,33M,0;
SRR038985.300      4          *          0           0          *          *          0
0          AAACTGCGGGGTCTCACTTTTTTGGGTTTGGGGT            124,/08/5&6-&,(;/4+
%7,+5.:1',*;8:&
```

61

Unmapped read

# BAM File Format (Alignment Format)

- BAM: BZGF compressed SAM Format

- Not human readable

- ~ 1 / 5 size of SAM

# gtf / gff3 files



| Chr1 | amel_OGSv3.1 | gene | 204921 | 223005 | . | + | . | ID=GB42165 |
| Chr1 | amel_OGSv3.1 | mRNA | 204921 | 223005 | . | + | . | ID=GB42165-RA;Parent=GB42165 |
| Chr1 | amel_OGSv3.1 | 3'UTR | 222859 | 223005 | . | + | . | Parent=GB42165-RA |
| Chr1 | amel_OGSv3.1 | exon | 204921 | 205070 | . | + | . | Parent=GB42165-RA |
| Chr1 | amel_OGSv3.1 | exon | 222772 | 223005 | . | + | . | Parent=GB42165-RA |

Chromosome ID · Source · Gene feature · Start location · End location · Score (user defined) · Strand · Phase · Attributes (hierarchy)

- Be aware of the format being used and its compatibility!
- Some tools will only work with gtf / gff3
- Ensembl GTF is NOT the same as UCSC GTF (even for the same assembly)

# GTF / GFF3 Fields

- https://www.ensembl.org/info/website/upload/gff.html

## Fields

Fields **must** be tab-separated. Also, all but the final field in each feature line must contain a value; "empty" columns should be denoted with a '.'

1. **seqname** - name of the chromosome or scaffold; chromosome names can be given with or without the 'chr' prefix. **Important note**: the seqname must be one used within Ensembl, i.e. a standard chromosome name or an Ensembl identifier such as a scaffold ID, without any additional content such as species or assembly. See the example GFF output below.

2. **source** - name of the program that generated this feature, or the data source (database or project name)

3. **feature** - feature type name, e.g. Gene, Variation, Similarity

4. **start** - Start position of the feature, with sequence numbering starting at 1.

5. **end** - End position of the feature, with sequence numbering starting at 1.

6. **score** - A floating point value.

7. **strand** - defined as + (forward) or - (reverse).

8. **frame** - One of '0', '1' or '2'. '0' indicates that the first base of the feature is the first base of a codon, '1' that the second base is the first base of a codon, and so on..

9. **attribute** - A semicolon-separated list of tag-value pairs, providing additional information about each feature.

# Summary

| Raw Data | | Trimming | | Mapping (index) | | Read Counting (gtf) |
|---|---|---|---|---|---|---|
| • Fastq | → | • Fastq | → | • sam / bam | → | • txt / csv |