

# RNA-Seq (2)

## BCB 5250 Introduction to Bioinformatics II

Spring 2020

Tae-Hyuk (Ted) Ahn

Department of Computer Science  
Program of Bioinformatics and Computational Biology  
Saint Louis University



SAINT LOUIS  
UNIVERSITY™

— EST. 1818 —

# Tools:

- Data Quality Control
- Read Mapping
- Differential Gene Expression

# Data Quality Control

## Quality assessment

- FastQC, MultiQC

## Trim and filtering:

- FASTX Tool Kit: [http://hannonlab.cshl.edu/fastx\\_toolkit/](http://hannonlab.cshl.edu/fastx_toolkit/)
- Trimmomatic: <http://www.usadellab.org/cms/?page=trimmomatic>
- Tim Galore: [http://www.bioinformatics.babraham.ac.uk/projects/trim\\_galore/](http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/)
- fastp: <https://github.com/OpenGene/fastp>
- BBTools – BBduk: <https://jgi.doe.gov/data-and-tools/bbtools/bb-tools-user-guide/bbduk-guide/>
- Cutadapt: <https://cutadapt.readthedocs.io/en/stable/>

## Error correction:

- BBTools – Tadpole: <https://jgi.doe.gov/data-and-tools/bbtools/bb-tools-user-guide/tadpole-guide/>
- <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-017-1784-8>
- For long sequencing, check the recently published papers and archive such as [http://www.pacb.com/asset\\_tags/error-correction/](http://www.pacb.com/asset_tags/error-correction/)

<https://galaxyproject.org/>

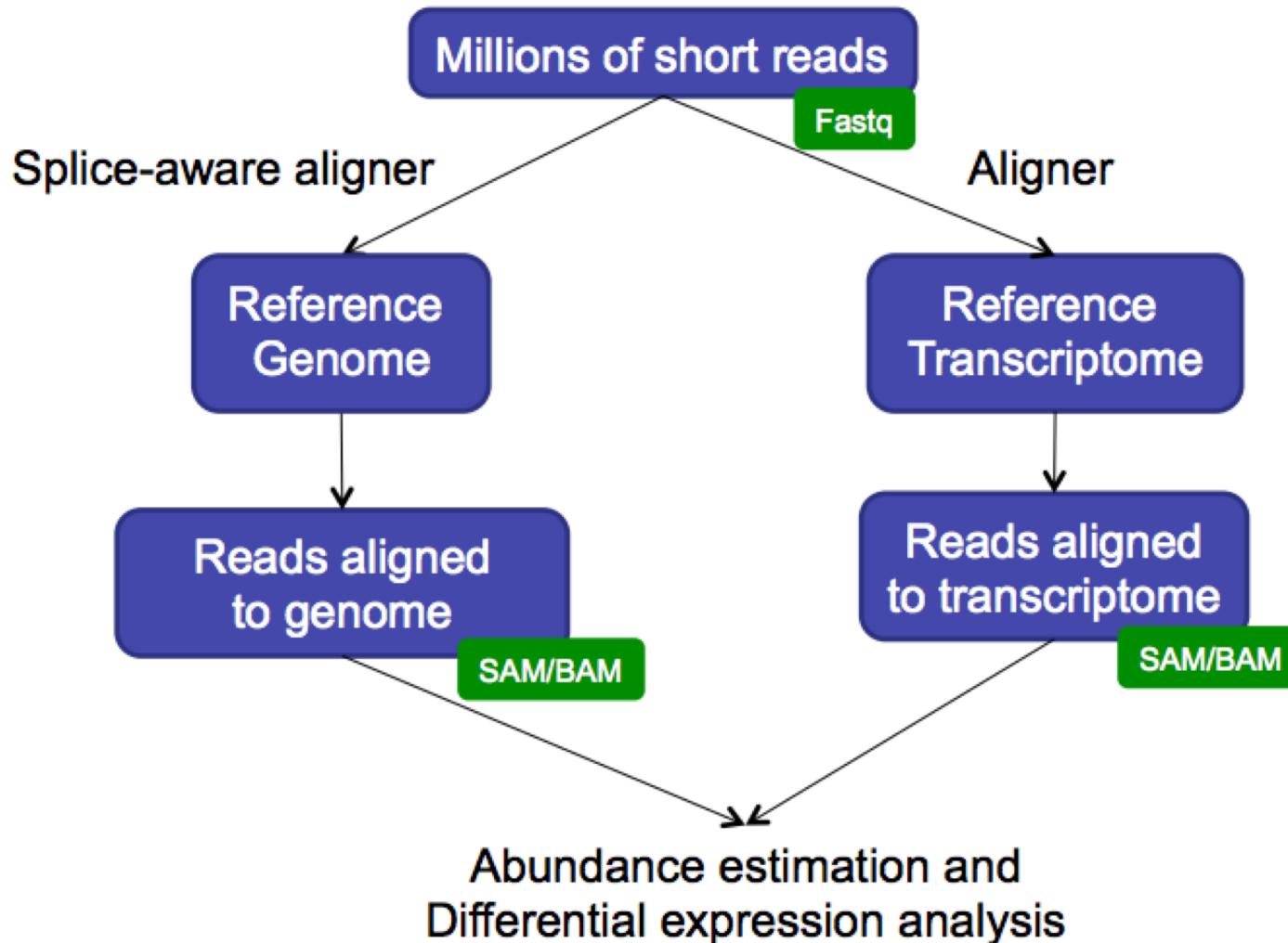
## Galaxy Community Hub

Galaxy is an open, web-based platform for accessible, reproducible, and transparent computational biomedical research.

- **Accessible:** Users without programming experience can easily specify parameters and run tools and workflows.
- **Reproducible:** Galaxy captures information so that any user can repeat and understand a complete computational analysis.
- **Transparent:** Users share and publish analyses via the web and create Pages, interactive, web-based documents that describe a complete analysis.

Welcome to the Galaxy Community Hub, where you'll find community curated documentation of all things Galaxy.

# Mapping Reads



# Aligner

- Short Read Mapper
  - Bowtie2: <http://bowtie-bio.sourceforge.net/bowtie2/index.shtml>
  - BWA: <http://bio-bwa.sourceforge.net/>
- Slice Aligner
  - TopHat2: <https://ccb.jhu.edu/software/tophat/index.shtml>
  - HiSAT2: <https://daehwankimlab.github.io/hisat2/>
  - STAR: <https://github.com/alexdobin/STAR>

# Mapping

- Input
  - Fastq files
  - Index of genome/transcriptome
  - Annotation file (optional for some, but required for others)
- Output
  - SAM (text) / BAM (binary) alignment files
    - SAMtools – SAM/BAM file manipulation (<http://samtools.sourceforge.net/>) (<http://www.htslib.org/>)
    - Picard-tools – SAM/BAM file manipulation (<https://broadinstitute.github.io/picard/>)
  - Summary statistics (per read library)
    - reads with unique alignment
    - reads with multiple alignments
    - reads with no alignment
    - reads properly paired (for paired-end libraries)

# Most Aligner

Most aligners have modules for below two steps:

1. Generating genome indexes files.
2. Mapping reads to the genome.

Why?

# Let us test bowtie2 aligner

1. Open manual (<http://bowtie-bio.sourceforge.net/bowtie2/manual.shtml#getting-started-with-bowtie-2-lambda-phage-example>)
2. Let us work “Getting Started”
  - Copy the genome to your lab directory and index a reference genome
    - /public/ahnt/courses/bcb5250/rna\_seq\_lab/bowtie2\_example/reference/lambda\_virus.fa
  - Aligning example reads
    - /public/ahnt/courses/bcb5250/rna\_seq\_lab/bowtie2\_example/reads/reads\_1.fq
  - Paired-end example
    - reads\_1.fq, reads\_2.fq
  - From the paired-end aligning output, let us test SAMtools/BCFtools downstream analysis
    - Convert SAM to BAM
    - Sort
    - To generate variant calls in VCF format
    - Then to view the variants

# SAM (BAM) Format

- Sequence Alignment/Map format
  - Universal standard
  - Human-readable (SAM) and compact & binary (BAM) forms
- Structure
  - Header
    - version, sort order, reference sequences, read groups, program/processing history
  - Alignment records

# Check header

```
$ samtools view -H -S XXX.sam
```

```
[samopen] SAM header is present: 1 sequences.
```

```
@HD VN:1.0SO:unsorted
```

```
@SQ SN:genome LN:451226
```

```
@PG ID:bowtie2 PN:bowtie2 VN:2.3.4.3 CL:"/usr/local/miniconda/bin/bowtie2-align-s --wrapper basic-0 -x ../GENOME_data/genome -S Sp_ds.sam -p 12 -t -1 Sp_ds.left.fq.gz -2 Sp_ds.right.fq.gz"
```

# Check alignment

```
$ samtools view -S XXX.sam | more
```

# Samtools

- <http://www.htslib.org/>
- <https://samtools.github.io/hts-specs/SAMv1.pdf>

← → C ⌂ ⓘ https://software.broadinstitute.org/software/igv/



The screenshot shows the IGV software interface. On the left is a sidebar with a logo, navigation links (Home, Downloads, Documents, etc.), and a search bar. The main area features a large title "Integrative Genomics Viewer" and a detailed genomic visualization with multiple tracks and data layers.

## Home

# Integrative Genomics Viewer



The main page has a "Home" section with the IGV logo and a "Search website" input field. Below it is an "Overview" section with a brief description of IGV and its availability in different forms. To the right is a "Citing IGV" section with citation information and a link to a PMC article.

**Overview**

The **Integrative Genomics Viewer (IGV)** is a high-performance visualization tool for interactive exploration of large, integrated genomic datasets. It supports a wide variety of data types, including array-based and next-generation sequence data, and genomic annotations.

IGV is available in multiple forms, including:

- the original **IGV** - a Java desktop application,
- IGV-Web** - a web application,
- IGV.js** - a JavaScript component that can be embedded

**Citing IGV**

To cite your use of IGV in your publication, please reference one or more of:

James T. Robinson, Helga Thorvaldsdóttir, Wendy Winckler, Mitchell Guttman, Eric S. Lander, Gad Getz, Jill P. Mesirov. [Integrative Genomics Viewer. Nature Biotechnology 29, 24–26 \(2011\). \(Free PMC article here\)](#).

Helga Thorvaldsdóttir, James T. Robinson, Jill P. Mesirov. [Integrative Genomics Viewer \(IGV\): high-performance genomics data visualization and exploration. Briefings in Bioinformatics 14, 470–490 \(2013\)](#)

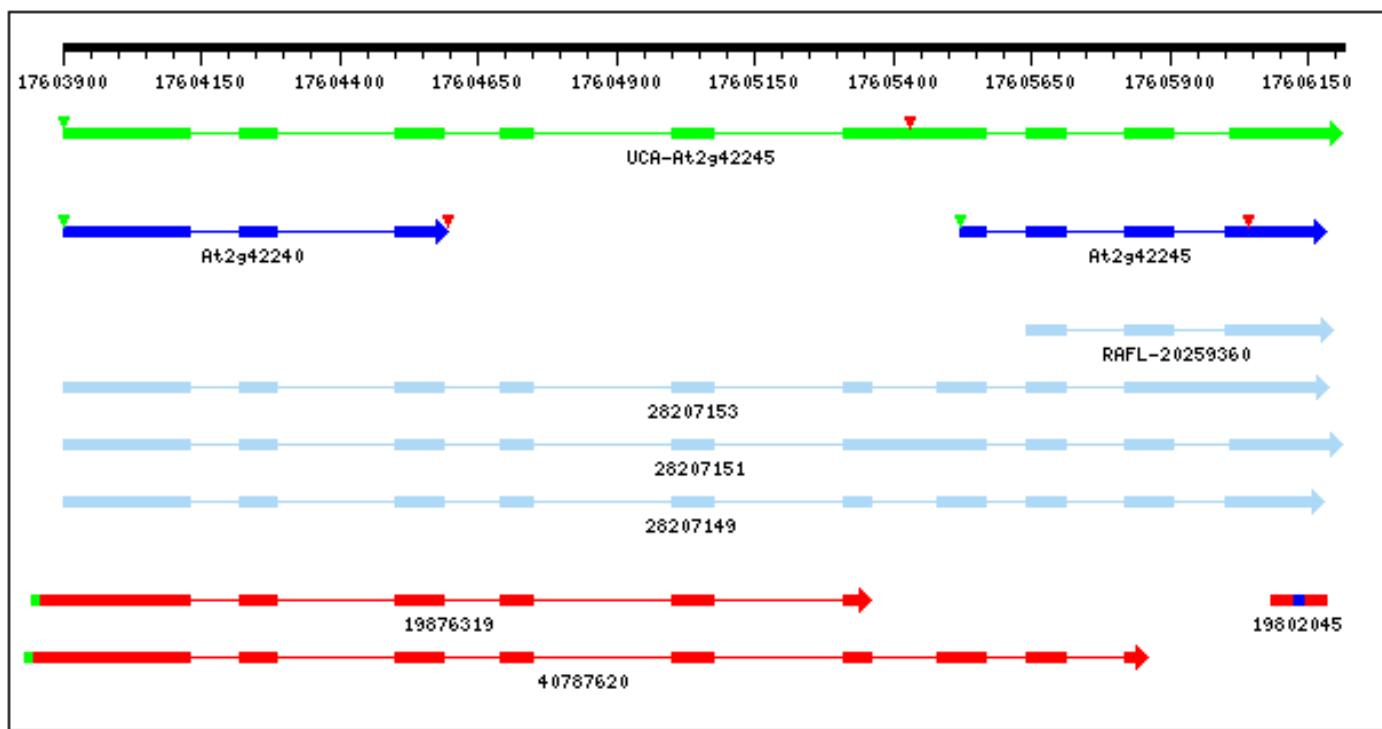
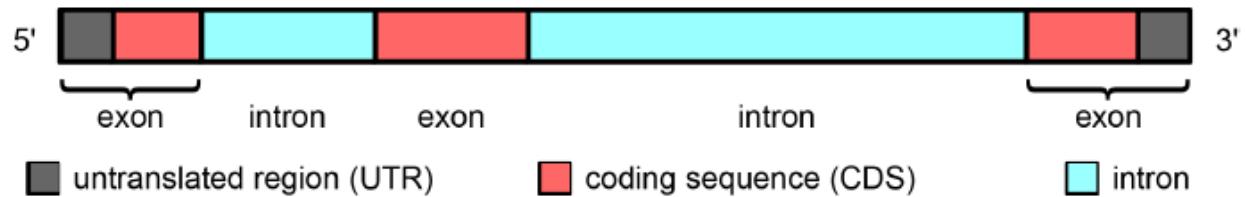
# Genome sequence file: where to get?

- Illumina iGenomes (recommended)
  - [http://support.illumina.com/sequencing/sequencing\\_software/igenome.html](http://support.illumina.com/sequencing/sequencing_software/igenome.html)
- Ensembl
  - <http://ensemblgenomes.org/info/access/ftp>
- NCBI genome
  - <http://www.ncbi.nlm.nih.gov/genome/>
- Organism specific databases/websites.

e!Ensembl



# Genome annotation file



# GFF/GTF File Format - Definition

- The GFF (General Feature Format) format
  - one line per feature
  - each containing 9 columns of data
  - plus optional track definition lines.
- GFF has many versions (GFF, GFF2, GFF3)
- GTF (General Transfer Format) identical to GFF2.
- Most spliced aligner supports both GTF and GFF3 (mostly GTF)

# GTF/GTF2 format

9 columns:

```
<seqname> <source> <feature> <start> <end> <score> <strand> <frame> [attributes]
```

- seqname - name of the chromosome or scaffold
- source – program or database that generated this feature.
- feature – Examples: "CDS", "gene", "transcript", and "exon".
- start - The starting position of the feature in the sequence.
- end - The ending position of the feature (inclusive).
- score - A score between 0 and 1000.
- strand – '+' (forward) or '-' (reverse) or '.' (don't know/don't care).
- Frame – reading frame '0', '1' or '2'
- attribute – A semicolon-separated list of tag-value pairs, providing additional information about each feature.

# Example of GTF2 format

```
AB000381 Twinscan CDS      380    401    .    +    0    gene_id "001"; transcript_id "001.1";
AB000381 Twinscan CDS      501    650    .    +    2    gene_id "001"; transcript_id "001.1";
AB000381 Twinscan CDS      700    707    .    +    2    gene_id "001"; transcript_id "001.1";
AB000381 Twinscan start_codon 380    382    .    +    0    gene_id "001"; transcript_id "001.1";
AB000381 Twinscan stop_codon 708    710    .    +    0    gene_id "001"; transcript_id "001.1";
```

A simple example with 3 translated exons. Order of rows is not important.

Some annotation sources (e.g. Ensembl) add the gene\_name attribute

```
gene_id "ENSBTAG00000020601"; transcript_id "ENSBTAT00000027448"; gene_name "ZNF366";
```

<http://mblab.wustl.edu/GTF2.html>

# Generic Feature Format Version 3 (GFF3)

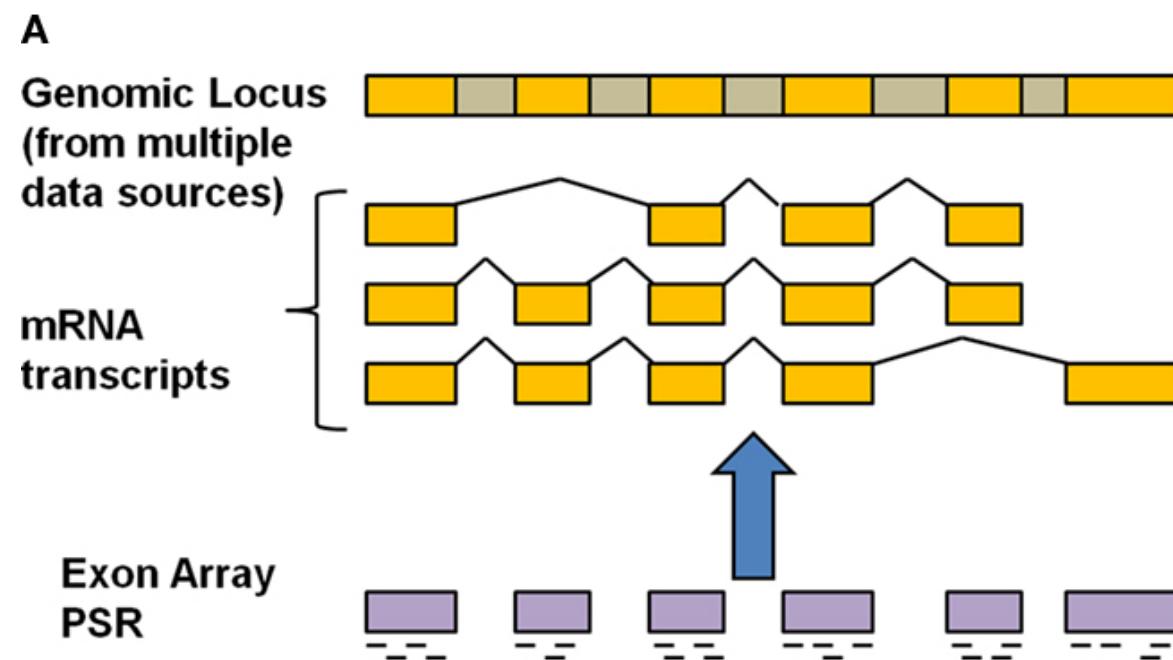
```
##gff-version 3.2.1
##sequence-region ctg123 1 1497228
ctg123 . gene    1000 9000 . + . ID=gene00001;Name=EDEN
ctg123 . TF_binding_site 1000 1012 . + . ID=tfbs00001;Parent=gene00001
ctg123 . mRNA    1050 9000 . + . ID=mRNA00001;Parent=gene00001;Name=EDEN
ctg123 . mRNA    1050 9000 . + . ID=mRNA00002;Parent=gene00001;Name=EDEN
ctg123 . mRNA    1300 9000 . + . ID=mRNA00003;Parent=gene00001;Name=EDEN
ctg123 . exon    1300 1500 . + . ID=exon00001;Parent=mRNA00003
ctg123 . exon    1050 1500 . + . ID=exon00002;Parent=mRNA00001,mRNA00002
ctg123 . exon    3000 3902 . + . ID=exon00003;Parent=mRNA00001,mRNA00002
ctg123 . exon    5000 5500 . + . ID=exon00004;Parent=mRNA00001,mRNA00002
ctg123 . exon    7000 9000 . + . ID=exon00005;Parent=mRNA00001,mRNA00002
ctg123 . CDS     1201 1500 . + 0 ID=cds00001;Parent=mRNA00001;Name=edEN
ctg123 . CDS     3000 3902 . + 0 ID=cds00001;Parent=mRNA00001;Name=edEN
ctg123 . CDS     5000 5500 . + 0 ID=cds00001;Parent=mRNA00001;Name=edEN
ctg123 . CDS     7000 7600 . + 0 ID=cds00001;Parent=mRNA00001;Name=edEN
ctg123 . CDS     1201 1500 . + 0 ID=cds00002;Parent=mRNA00002;Name=edEN
ctg123 . CDS     5000 5500 . + 0 ID=cds00002;Parent=mRNA00002;Name=edEN
ctg123 . CDS     7000 7600 . + 0 ID=cds00002;Parent=mRNA00002;Name=edEN
ctg123 . CDS     3301 3902 . + 0 ID=cds00003;Parent=mRNA00003;Name=edEN
ctg123 . CDS     5000 5500 . + 1 ID=cds00003;Parent=mRNA00003;Name=edEN
ctg123 . CDS     7000 7600 . + 1 ID=cds00003;Parent=mRNA00003;Name=edEN
ctg123 . CDS     3391 3902 . + 0 ID=cds00004;Parent=mRNA00003;Name=edEN
ctg123 . CDS     5000 5500 . + 1 ID=cds00004;Parent=mRNA00003;Name=edEN
ctg123 . CDS     7000 7600 . + 1 ID=cds00004;Parent=mRNA00003;Name=edEN
```

GFF3 adds parent feature

<http://www.sequenceontology.org/gff3.shtml>

# GFF2 vs GFF3

- GFF2
  - two-level hierarchies *transcript* → *exon*
- GFF3
  - three-level hierarchy of *gene* → *transcript* → *exon*



<http://gmod.org/wiki/GFF2>

# Conversion GFF3 To GTF

- Optional
- Use `gffread` (comes with the Cufflinks software suite)

```
$ gffread my.gff3 -T -o my.gtf
```

- See `gffread -h` for more information

# Download Prepare GTF/GFF file

- Download it from Illumina's iGenomes project (for model species)
  - [http://support.illumina.com/sequencing/sequencing\\_software/igenome.html](http://support.illumina.com/sequencing/sequencing_software/igenome.html)

Or

- Download gff3 files from genome database,
  - Ensembl gnome
  - <http://ensemblgenomes.org/>