

RNA-Seq: Aligner, Output (SAM), and Visualization

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 —

RNA-Seq Tools for:

- 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/
- Trimmomatic: <http://www.usadellab.org/cms/?page=trimmomatic>
- Tim 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/

<https://galaxyproject.org/>

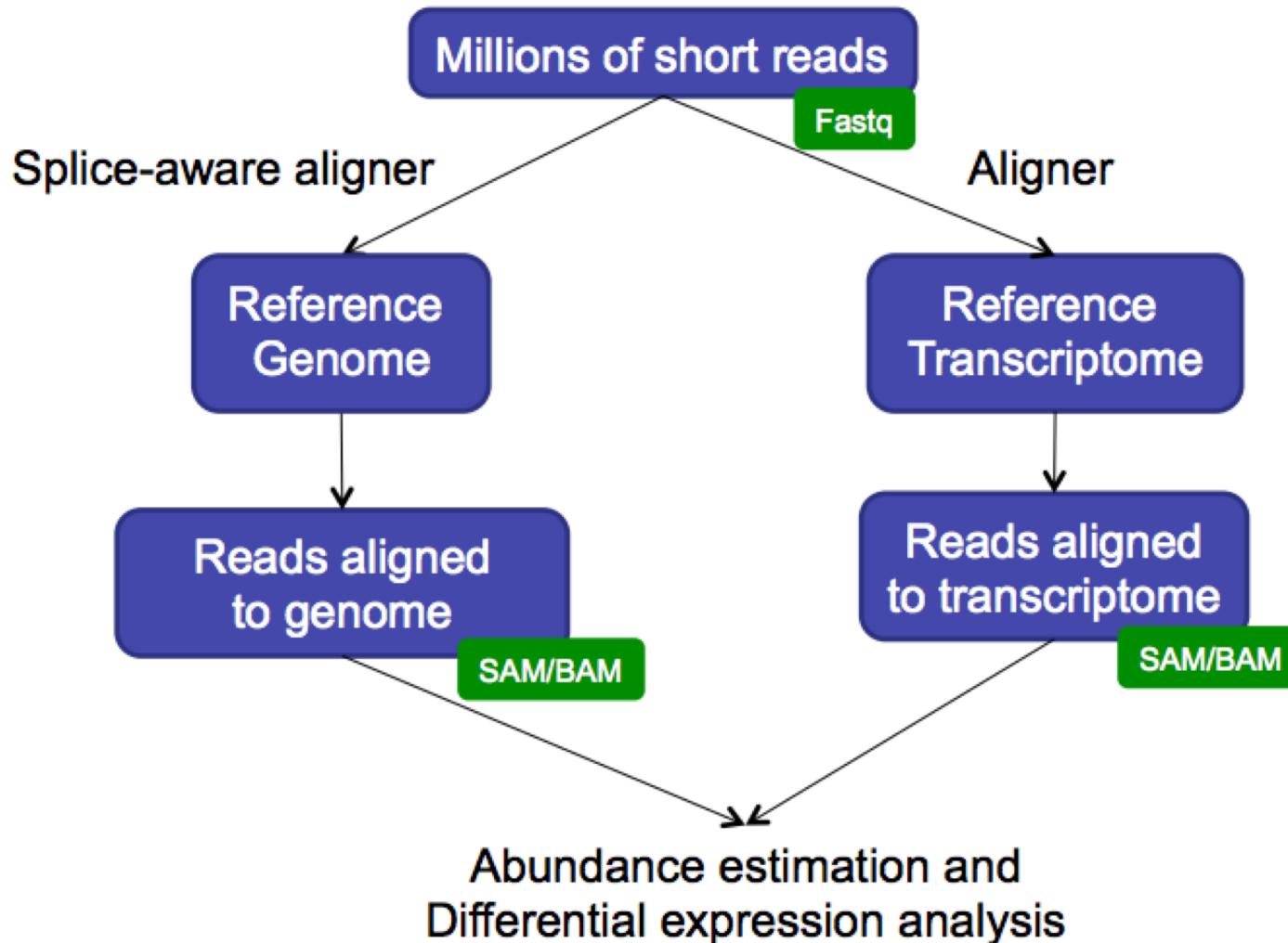
Galaxy Community Hub

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/>
- Splice 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/



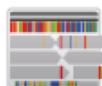
The screenshot shows the IGV software interface. On the left is a sidebar with a logo, navigation links (Home, Downloads, Documents, etc.), and a search bar. The main area features a large title "Integrative Genomics Viewer" and a detailed genomic visualization with multiple tracks and data layers.

Home

Integrative Genomics Viewer



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 \(2013\)](#)

- Load Genome
- Load sorted BAM alignments (you need index file)