

Nucleotide sequence analysis

Bioinformatics, 2nd year
UP FAMNIT

Plan for lectures

- RNA-Seq (transcriptome analysis)
- Genome assembly
- SNP calling
- Metabarcoding
- Genome wide association analysis

Literature

- Pevsner, J. 2015. Bioinformatics and functional genomics. Wiley Blackwell.
- Material for practicals (e-classroom)
- Korpelainen et. al. 2015. RNA-seq Data Analysis: A Practical Approach.
- István Albert. The Biostar Handbooks
<https://www.biostarhandbook.com/>
- scientific articles

ishán Albert, Bioinformatics, Penn State

Home

Publications

Software

Consulting

Bioinformatics education is severely lacking across the world. Most course work is inadequate: practitioners, inadequately trained, propagating misconceptions and fallacies. Substantial resources are spent on ineffective training: needless travel, expensive week-long workshops, "bootcamps for summer" type of training, promoting "click-the-button" kind of solutions unceremoniously, with little effect.

After waiting a decade for a solution, I have decided to write the materials that I felt were missing. This is the genesis of the **Biostar Handbooks**, initially a single book, now a series of volumes focusing solely on the practical aspects of bioinformatics data analysis.

The Biostar Handbooks

The **Biostar Handbooks** introduce readers to **bioinformatics**, the scientific discipline at the intersection of biology, computer science, and statistical data analysis dedicated to the digital processing of genomic information.

- The **Biostar Handbook** - An introduction to bioinformatics as a scientific field.
- The **Art of Bioinformatics Scripting** - Learn advanced Linux and Bash scripting skills.
- **How to Use Bioinformatics** - Master data analysis.
- **Genomic Data Science Analysis** - Advanced topics devoted to the study of the Genomes.
- **Biostar Workflows** - Best practices for developing bioinformatics pipelines.

The **Biostar Handbooks** deliver simple, concise, and relevant information for those looking to understand the field of bioinformatics as a data science. It is a comprehensive, practical handbook that covers all major application areas of bioinformatics.

Teaching

I teach the **BBB 402: Applied Bioinformatics** course as resident instructor in the fall.

Dr. Ishán Albert
Assistant Professor of Bioinformatics
Department of Biochemistry and Molecular Biology
Pennsylvania State University

3500 Life Sciences
PA 16802-3000
USA
ishan@psu.edu

<https://www.biostarhandbook.com/>

Obligations and Grading system

Please be advised that attendance at practical sessions is mandatory.
The final grade will be composed based on the following criteria:
70% Project work
30% Oral exam

Student's obligations:
Each student must analyze their own dataset and prepare a laboratory notebook with a documented workflow.
Include all commands and explain them, including options (flags) selected with each program. Document your work with screenshots and output files. Interpret the results, including graphs obtained with FastQC. Report the number of reads in your dataset, how many were removed after quality filtering, etc. Answer all questions from the worksheets.
Laboratory notebooks should be uploaded to the e-classroom.
The report must be submitted at least 1 week before the exam date.
Please note: You cannot attend the oral exam if you haven't submitted the notebook before the deadline!
Oral Defense:
During the oral defense, expect questions related to the project content or general questions about nucleotide sequence analysis, such as "Describe the SAM file! How can it be converted to a BAM file?" and questions related to the lectures.
The oral defense will occur on the date specified in SIS.

Transcriptome analysis (RNA-Seq)

2

Low- and high-throughput technologies to study mRNAs

Three techniques for the study of mRNA:

- complementary DNA (cDNA) libraries
- microarrays (e.g. using the Affymetrix platform)
- RNA-seq (Chapter 11)

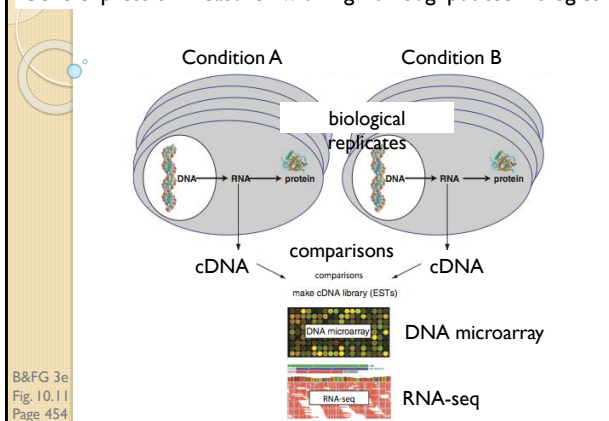
Low throughput techniques (Northern blots, PCR) may seem laborious and able to provide only limited amounts of information.

Yet they also serve as trusted “gold standards” and provide crucial validation of high throughput techniques.

B&FG 3e
Page 452

Source: Pevsner: Bioinformatics and Functional Genomics, 3rd Edition

Gene expression measured with high-throughput technologies

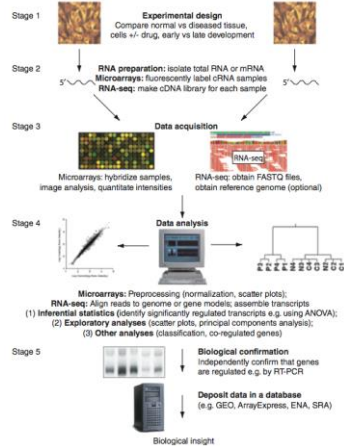


B&FG 3e
Fig. 10.11
Page 454

Transcriptome analysis (RNA-Seq)

- High throughput RNA sequencing (RNA-Seq) offers possibility to investigate the expression profiles at the transcriptional level and also identifying novel and non-coding transcripts
- The **reference-based** transcriptome analysis method is based on aligning the sequenced reads to a pre-existing reference genome, (followed by assembling overlapping alignments into transcripts).
- In contrast, the reference-free **de novo** transcriptome analysis method allows to directly assemble sequenced reads into transcripts by using high levels of redundancy and overlapping of reads, without using a reference genome.

Overview of the process of generating high-throughput gene expression data using microarrays or RNA-seq



B&FG 3e
Fig. 10.13
Page 462

Alignment based vs classification based RNA-Seq

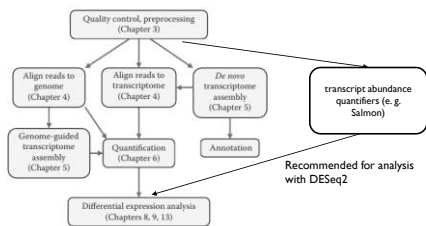


FIGURE 2.1 Possible paths in RNA-seq data analysis. In the beginning, the quality of reads is checked and if necessary, reads are preprocessed to remove low-quality data and artifacts. The read's origin is then identified by aligning them to a reference genome if available. Novel genes and transcripts are detected using genome-guided transcriptome assembly, and gene and transcript expression are quantified. Alternatively, gene and transcript discovery can be skipped and expression quantified only for known genes and transcripts. If the reference genome is not available, reads can be aligned and quantified using a reference transcriptome instead. If a transcriptome is not available, it can be produced from reads using *de novo* transcriptome assembly. When abundance estimates are obtained using one of these paths, expression differences between sample groups can be analyzed using statistical testing. The details of each step can be found in the chapters indicated in parentheses.

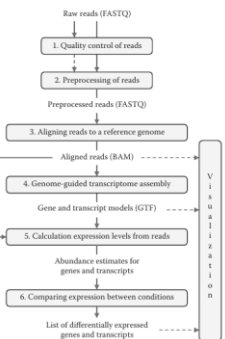


FIGURE 2.3 Differential expression analysis workflow consists of several, inter-related steps. The typical output file formats are indicated in parentheses.

What are the tradeoffs?

- Genome base methods:
1. Allow us to visualize the data in the context of the entire genome.
 2. Enable us to discover/validate new transcripts.

- Transcriptome and classification based methods:
1. Are typically more accurate.
 2. Require lower computational resources.

Counting Reads per Genes

- Counting Reads per Genes
- The simplest way of estimating expression is to count reads per genes. Several tools are available: HTSeq, Cufflinks, RSubread (Bioconductor package)
- Counting Reads per Transcripts
- Counting Reads per Exons

What is the final result of an RNA-Seq analysis?

The result of an RNA-Seq analysis is a *quantification matrix*. For our toy example, the file might look like this:

name	control	shock
Gene A	100	200
Gene B	80	60
Gene C	120	180
...		

[illegible]

SRA toolkit:

fastq-dump to obtain FASTA formatted data

```
$ fastq-dump -X 3 -Z SRR390728 -fasta 36
Read 3 spots for SRR390728
Written 3 spots for SRR390728
>SRR390728.1 1 length=72
CATTCTTCACGTAGTCTCTGAGCCCTGGTTTTTCAGC
GATGGAGAATGACTTTGACAAGCTGAGAGAAGNTNC
>SRR390728.2 2 length=72
AAGTAGGCTCTCGTCTGTGTTTTCTACGAGCTGTGT
TCCAGCTGACCCACTCCCTGGGTGGGGGGACTGGGT
>SRR390728.3 3 length=72
CCAGCTGGCCACAGAGTGTATCCCCGTTTTTACT
TATTTATTATTATTATTTTGGAGACAGACATTGGTC
```

B&FG 3e
Page 393

FASTQ format

The FASTQ format stores DNA sequence data as well as associated Phred quality scores of each base.

```
@EAS54_6_R1_2_1_413_324
CCCTTCTTGTCCTTCAGCGTTTCTCC
+
::3::::::::::::7:::::::::88
```

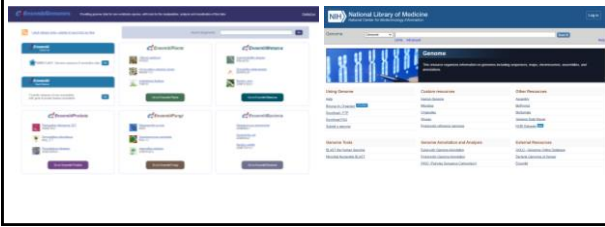
DNA read

Base quality score

```
@EAS54_6_R1_2_1_540_792
TTGGCAGGCCAAGGCCGATGGATCA
+
::::::::::::7::::-::3:83
@EAS54_6_R1_2_1_443_348
GTTGCTTCTGGCGTGGGTGGGGGGG
+EAS54_6_R1_2_1_443_348
::::::::::::9:7::,7:393333
```


Where genomes are available?

International Nucleotide Sequence Database Collaboration (INSDC) members [National Center for Biotechnology Information (NCBI), European Bioinformatics Institute (EBI), and DNA Data Bank of Japan (DDBJ)]



Where genomes are available?



<https://ngdc.cncb.ac.cn/gwh/>

Alignment based vs classification based RNA-Seq

- Aligners
- HISAT2
 - STAR
 - Bowtie (a splice unaware aligner)
 - BWA (a splice unaware aligner)
 - ...

A classification-based RNA-Seq works via so-called “pseudo-alignments”.

- Pseudo-alignment tools
- Kallisto
 - Salmon

Indexing reference genome/transcriptome

- Indexing reference genome
 - Incorporating the annotation into the index

Annotation file



The [formal specification of GFF3](#) is on the [Sequence Ontology](#) web site.



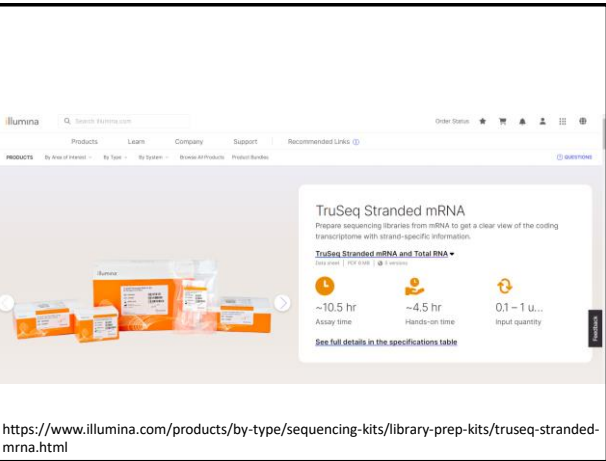
Stranded vs non-stranded RNA-Seq

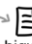

- Stranded (strand-specific) RNA-Seq

Preserves strand information—all mapped reads are aligned in the direction of transcription relative to the chromosomal strand

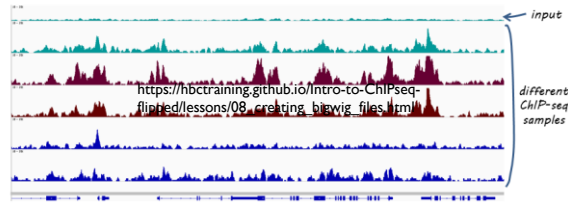
Ion Total RNA-Seq Kit v2

https://www.thermofisher.com/document-connect/document-connect.html?url=https://assets.thermofisher.com/TFS-Assets%2FLSG%2Fmanuals%2FMAN00010654_IonTotalRNASeqKit_v2_UG.pdf





for visualizing continuous data, e.g. in the UCSC Genome Browser or IGV, bigwig files come in really handy

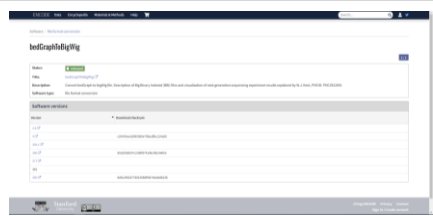


https://hbctraining.github.io/Intro-to-ChIPseq-flipped/lessons/08_creating_bigwig_files.html

remember that there are 2 deepTools for bam → bigwig conversion:

- ❖ `bamCoverage`: for individual files (like those shown here)
- ❖ `bamCompare`: to normalize two files to each other

https://hbctraining.github.io/Intro-to-ChIPseq-flipped/lessons/08_creating_bigwig_files.html

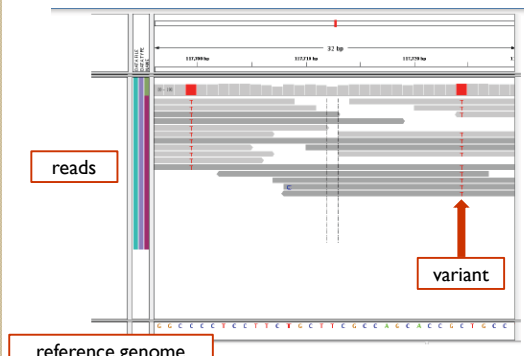


```
# Turn each BAM file into bedgraph coverage. The files will have the .bg extension
cat $(ls *.bam | parallel 'bedtools genomecov -ibam {} -o bamplit -bg > bam/{}.bg')
# Convert each bedgraph coverage into bigwig coverage. The files will have the .bw
cat $(ls *.bg | parallel 'bedgraphToBigwig bam/{} -i BEDD -ai bam/{} -o')
The resulting bedgraph and bigwig files will have *.bg and *.bw extensions
and are placed in the bam directory.
You can drag your bigwig file in the IGV* panels, and the coverage information
will load up much much faster. Below I have loaded all samples, rotated
the tracks, and colored them by samples. I have also turned on logarithmic
and automatic scaling. The resulting browser track is quite informative:
```

Biostar handbook

<https://www.encodeproject.org/software/bedgraphToBigwig/>

Next-generation sequence data: visualizing of short reads aligned to a reference genome



reads

variant

reference genome

BWA and other aligners produce output in the SAM format

Column	Description
1 QNAME	Query (pair) NAME
2 FLAG	bitwise FLAG
3 RNAME	Reference sequence NAME
4 POS	1-based leftmost POSition/coordinate of clipped sequence
5 MAPQ	MAPping Quality (Phred-scaled)
6 CIGAR	extended CIGAR string
7 MRNM	Mate Reference sequence Name ('=' if same as RNAME)
8 MPOS	1-based Mate POSition
9 ISIZE	Inferred insert SIZE
10 SEQ	query SEQUENCE on the same strand as the reference
11 QUAL	query QUALity (ASCII-33 gives the Phred base quality)
12 OPT	variable OPTional fields in the format TAG:VTYPE:VALU

<https://www.samformat.info/sam-format-flag>

<https://broadinstitute.github.io/picard/explain-flags.html>

Sequence alignment/map format (SAM) and BAM

- SAM is a common format having sequence reads and their alignment to a reference genome.
- BAM is the binary form of a SAM file.
- Aligned BAM files are available at repositories (Sequence Read Archive at NCBI, ENA at Ensembl)
- SAMTools is a software package commonly used to analyze SAM/BAM files.
- Visit <http://samtools.sourceforge.net/>

Anatomy of a Sequence Alignment/Map (SAM) file

- (1) The query name of the read is given (M01121...)
- (2) The flag value is 163 (this equals $1+2+32+128$)
- (3) The reference sequence name, chrM, refers to the mitochondrial genome
- (4) Position 480 is the left-most coordinate position of this read
- (5) The Phred-scaled mapping quality is 60 (an error rate of 1 in 10^6)
- (6) The CIGAR string (148M2S) shows 148 matches and 2 soft-clipped (unaligned) bases
- ```
home/bioinformatics$ samtools view 030c_s7.bam | less
M01121:5:000000000-AZDTN:1:2111:20172:15571 163 chrM
480 60 148M2S = 524 195 AATCTCATCAAT
ACAAACCCCTGCGCCATCTTACCAGACACACACACCGCTGCTAACCCCATACCCGAACC
AACCAAAACCCCAAAACACACCCCCACGCTTATGTAAGCTTACCTCTCTCAAGCAATACG
TGAAAAATGTTAGACGGG BBBBFBFB58FTGGGFGGGGGAACGHPHFGGAGFTH
AEFDGGT8?EGGGFGHGFHF?FFCHFH0E8EGFGGEE1FFEEH8GFFFGGGG8</0
1BG212222>F218F11FGFG1817GC<G111?FGDGGF=GHFFHC.-
RG:Z:Sample7 XC:1:148 XT:A:U NM:1:3 SM:1:37
AM:i:37 X0:i:1 X1:i:0 XM:i:3 XO:i:0 XG:i:0 MD:Z:19C109C0A17
```
- (7) An = sign shows that the mate reference matches the reference name
- (8) The 1-based left position is 524
- (9) The insert size is 195 bases
- (10) The sequence begins AATCT and ends ACGGG (its length is 150 bases)
- (11) Each base is assigned a quality score (from 88888 ending 8AC.-)
- (12) This read has additional, optional field at accompany the MiSeq analysis

(1) The query name of the read is given (MO1121...)

In this example we'll look at a file called `030c_s7.bam`. It is a **BAM** file (the binary of a **SAM**). Most software manipulates **BAM** files rather than **SAM**.

The `|` symbol (called “pipe”) indicates to send the results to another program—in this case to the utility called `less` that displays one page at a time on your terminal.

(10) The sequence begins AATCT and ends ACGGG (its length is 150 bases)

(11) Each base is assigned a quality score (from 000000 ending FHC. -)

additional, optional  
file to accompany  
the MiSeq analysis

B&FG 3e  
Fig.9-13  
Page 403

[illegible]

B&FG 3e  
Fig.9-14  
Page 404