

RNA-Sequencing Data Analysis

**Korea Institute of Science and Technology (KIST)
Clean Energy Research Center**

Sungmin Hwang

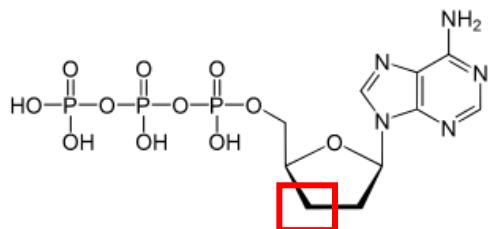
2022. 2. 11

DNA sequencing technology: 1st generation (Sanger's chain termination)

Original DNA sequence

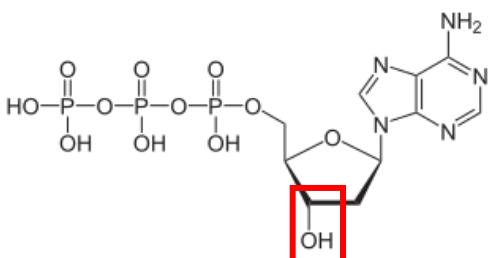
ATGCAGCGTTACCATG...

Amplification
+ Pol
+ ddNTPs

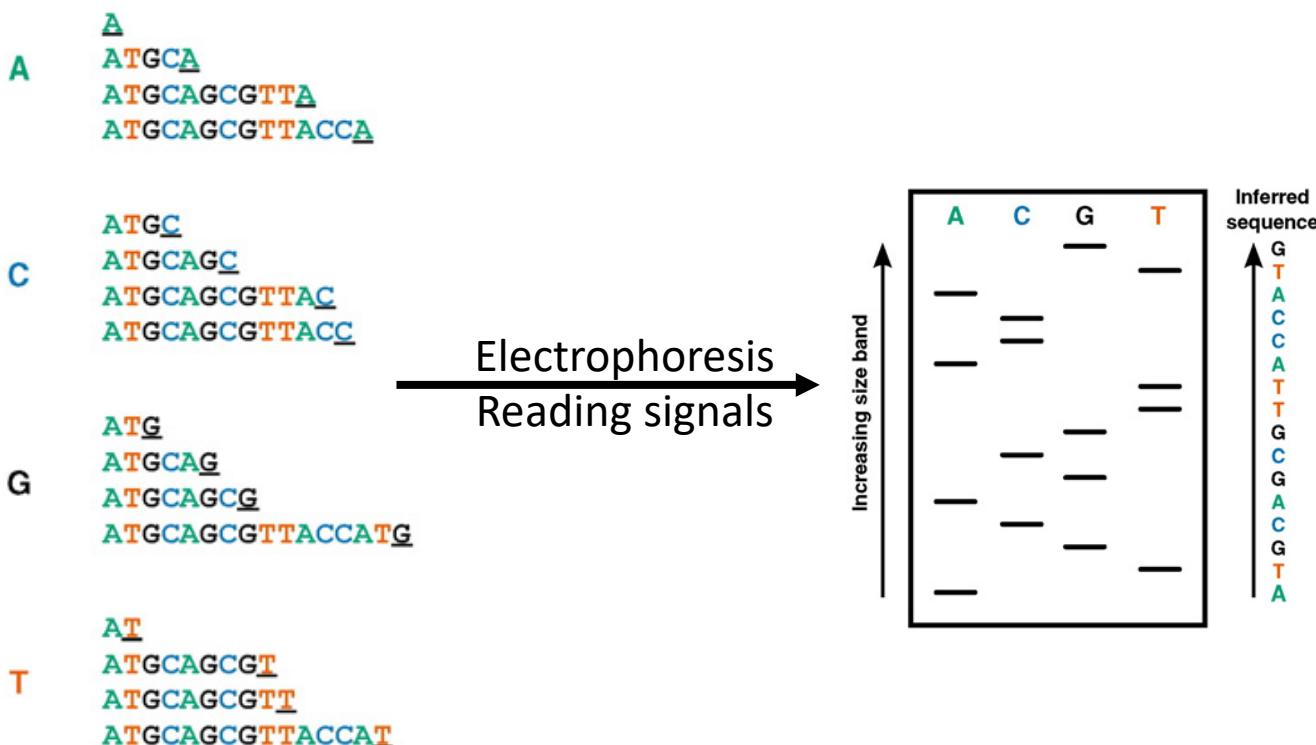


dideoxynucleotides (ddNTPs)

: lack the 3' hydroxyl group that is required for extension of DNA chains



deoxyribonucleotides (dNTPs)

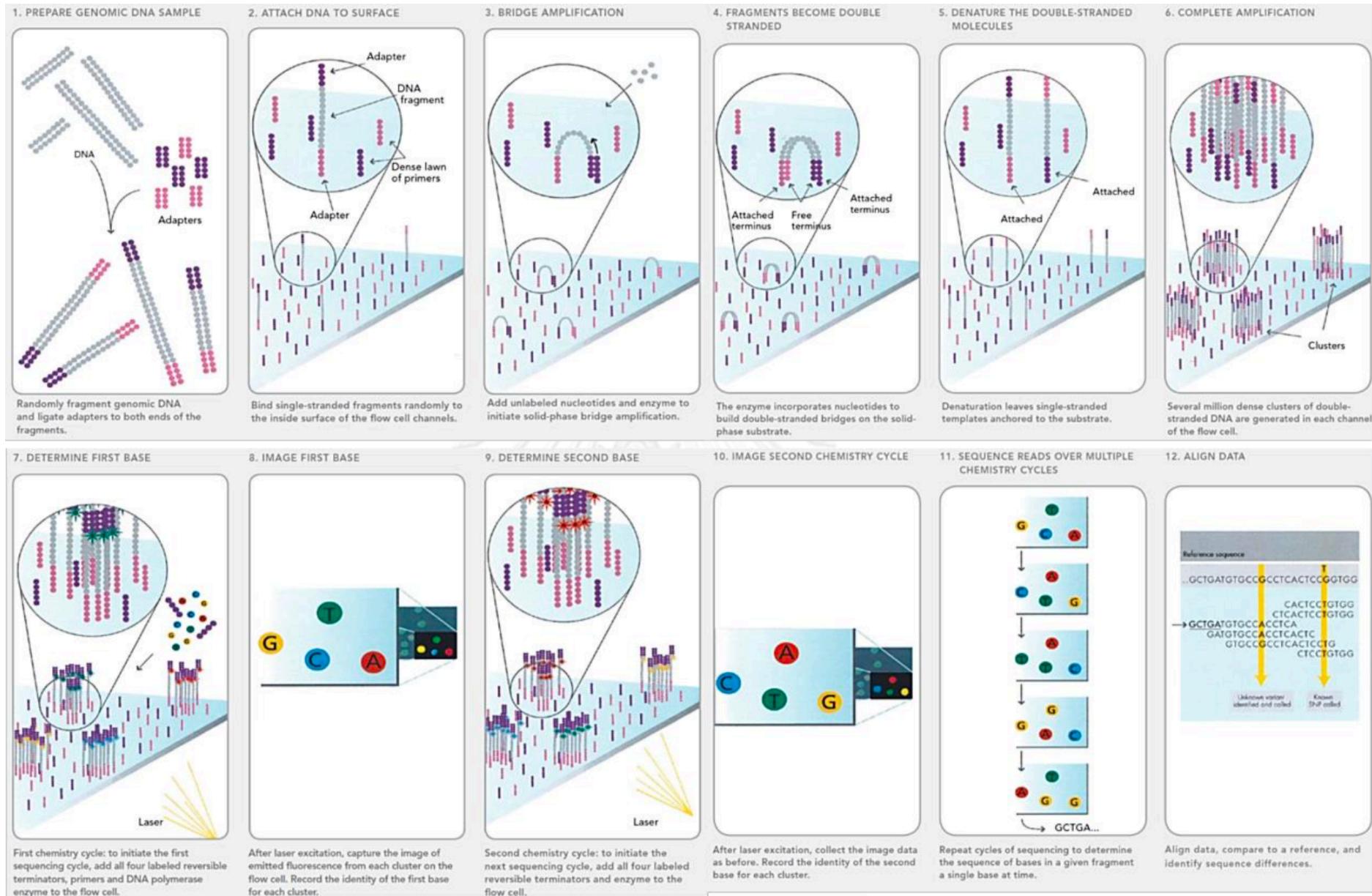


Limitation

- Coverage: ~ 1 kbp
- Organisms that have a big genome size
i.e., *E. coli* (4.6 Mbp), *Homo sapiens* (3,088 Mbp)

DNA sequencing technology: Next-generation (high-throughput)

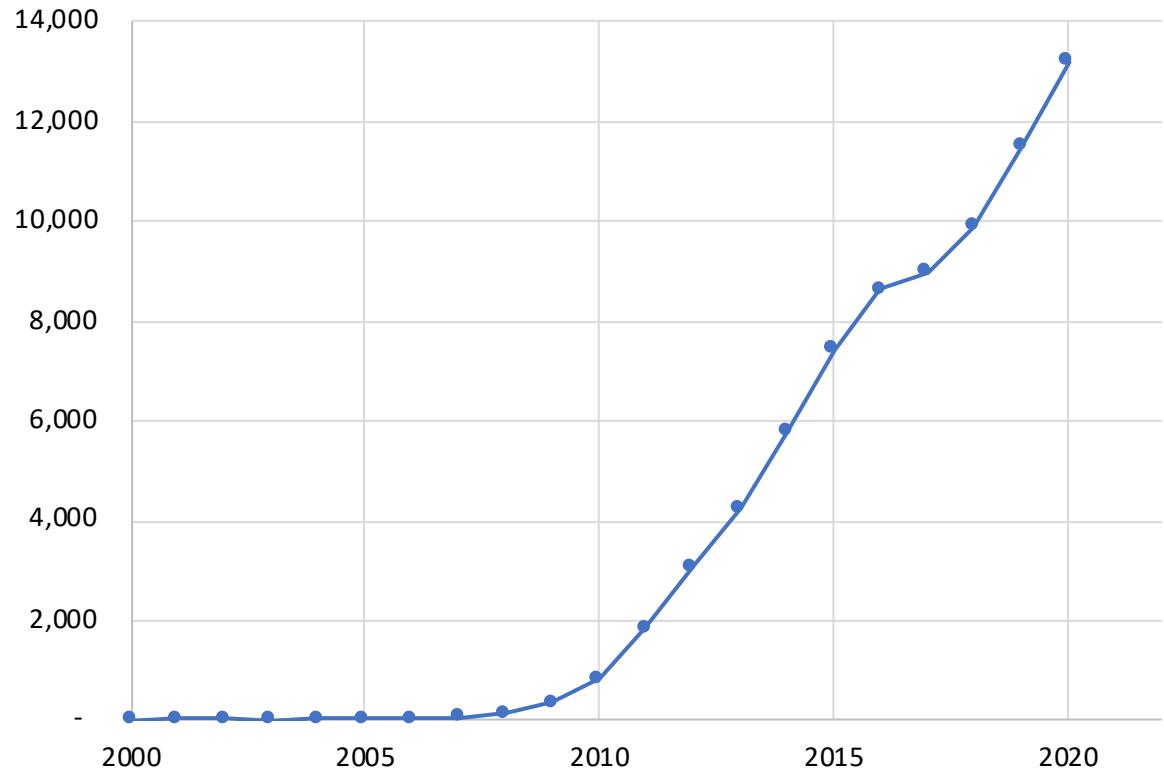
Illumina Co.,



Flow-cell

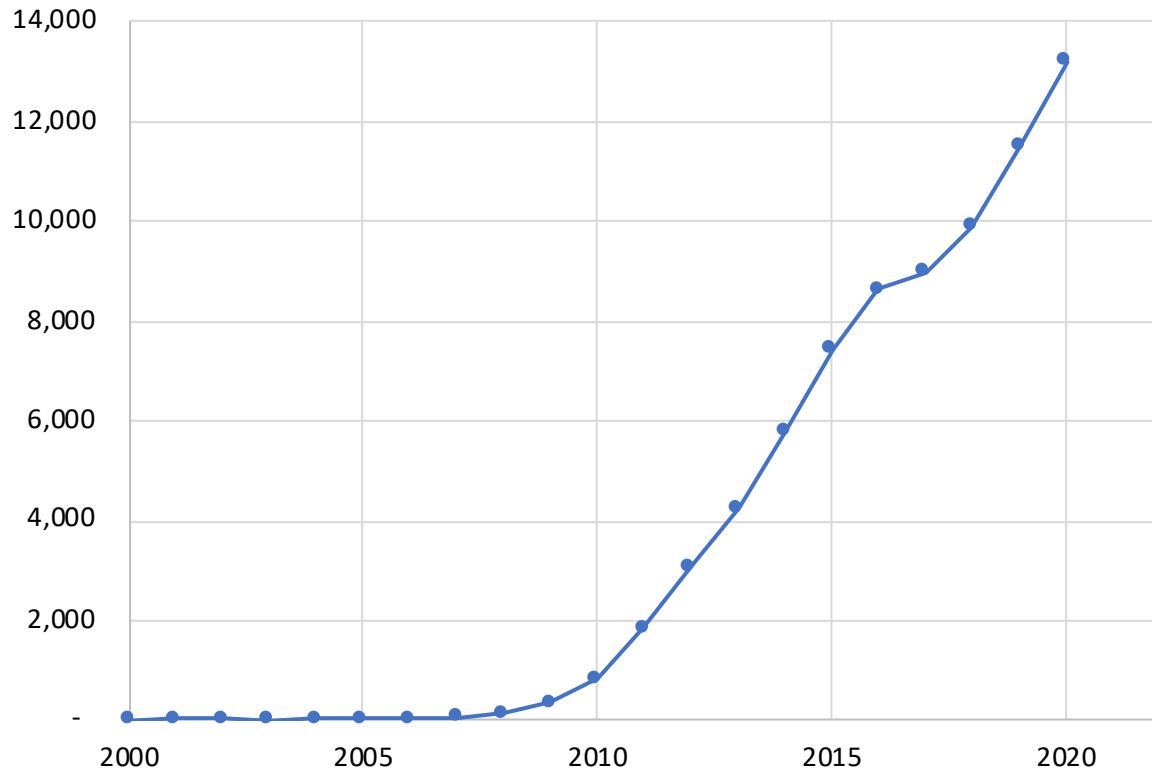
DNA sequencing technology: Next-generation (high-throughput)

No. of studies in Pubmed during 2000-2020



DNA sequencing technology: Next-generation (high-throughput)

No. of studies in Pubmed during 2000-2020



Illumina, Inc.

\$367.86 ↑1,778.75% +348.28 MAX

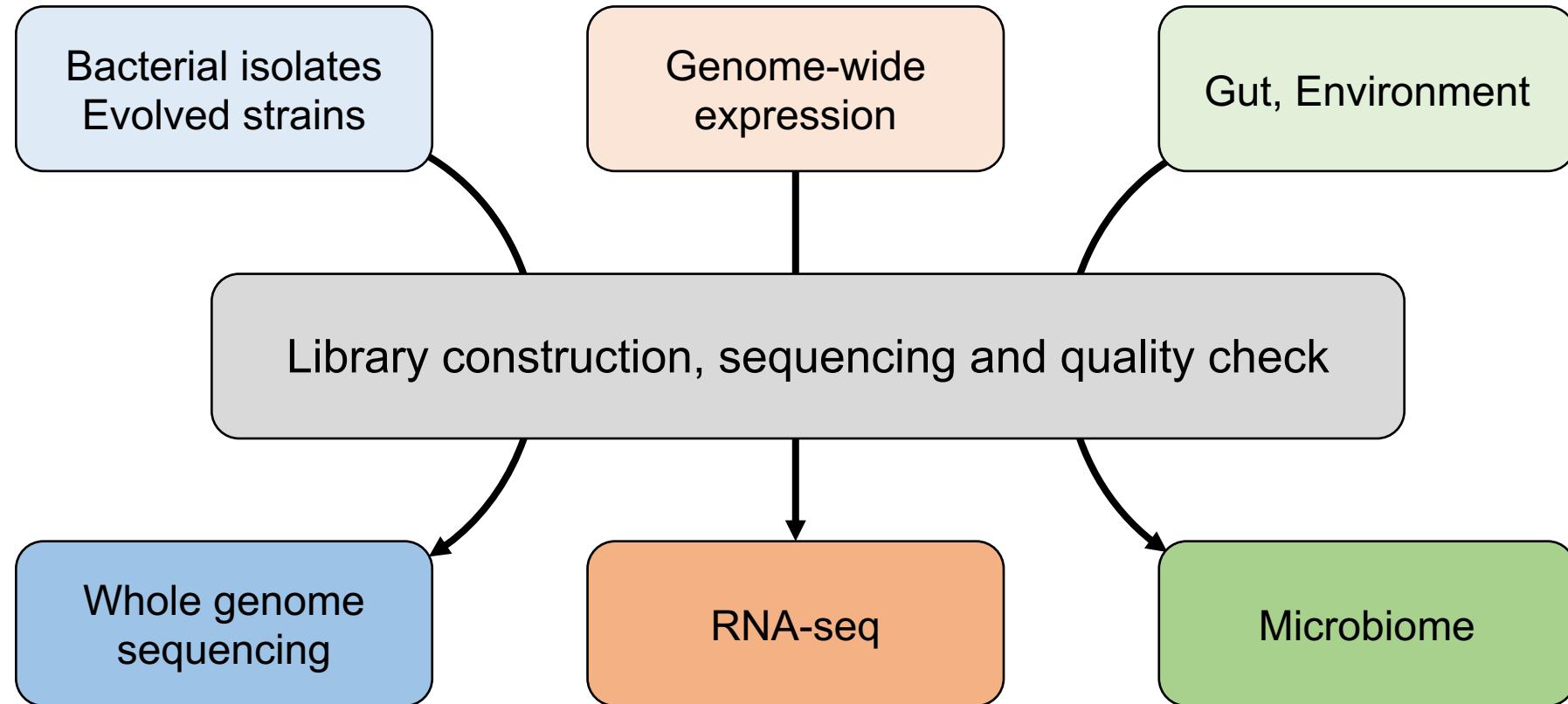
After Hours: \$367.86 (0.00%) 0.00

Closed: Feb 9, 5:22:07 PM UTC-5 · USD · NASDAQ · Disclaimer

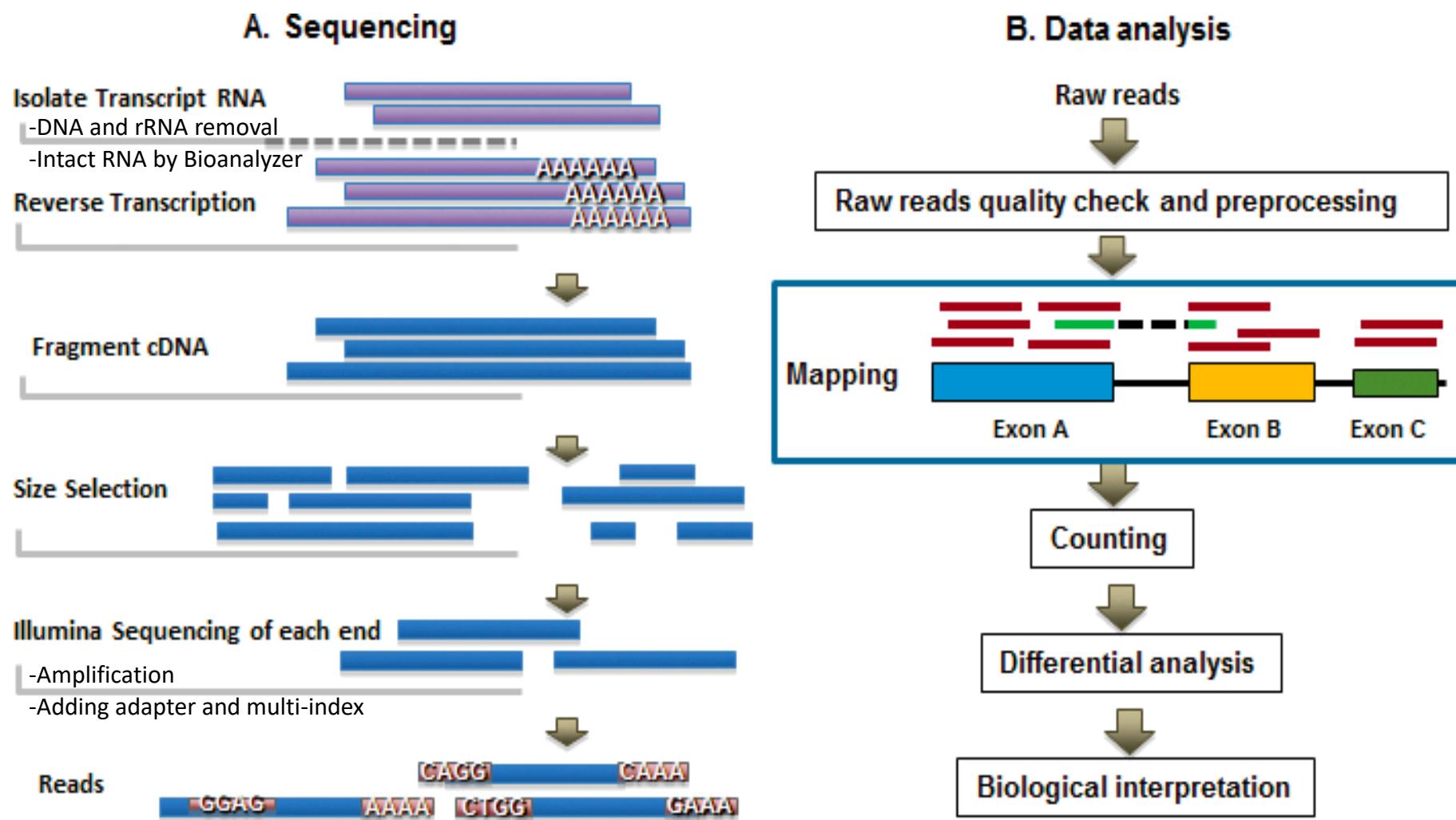


* 본 자료는 NGS 기술을 활용한 연구 내용 전달에 있으며,
투자 권유 및 종목 추천이 아님을 명시합니다.

Applications by using NGS



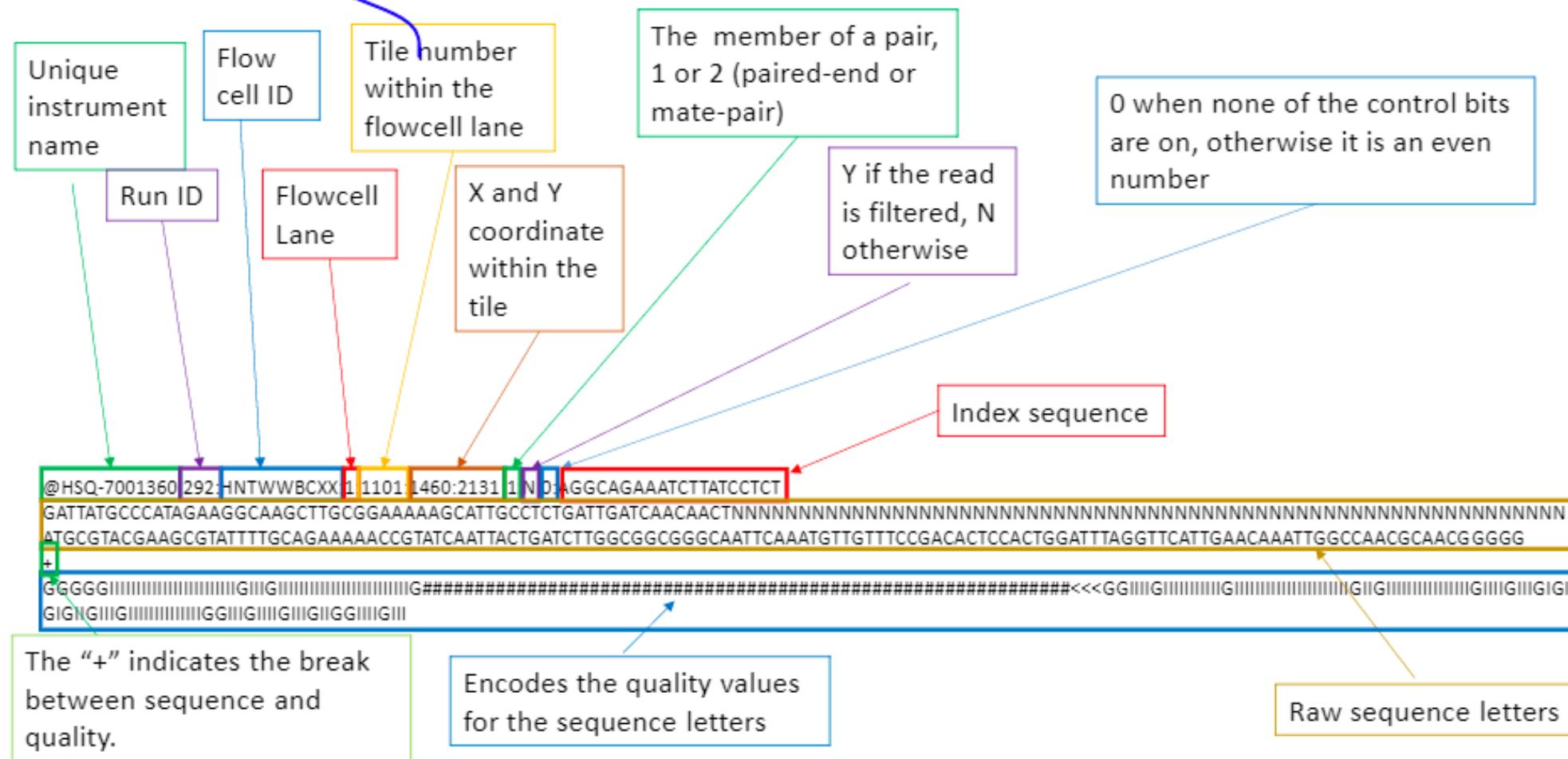
Overview of RNA-seq flowchart



Raw data format of NGS data

Raw reads → Quality check → Mapping → Counting → DEG analysis → Functional enrichment → Data interpretation

FASTQ File Format Analysis



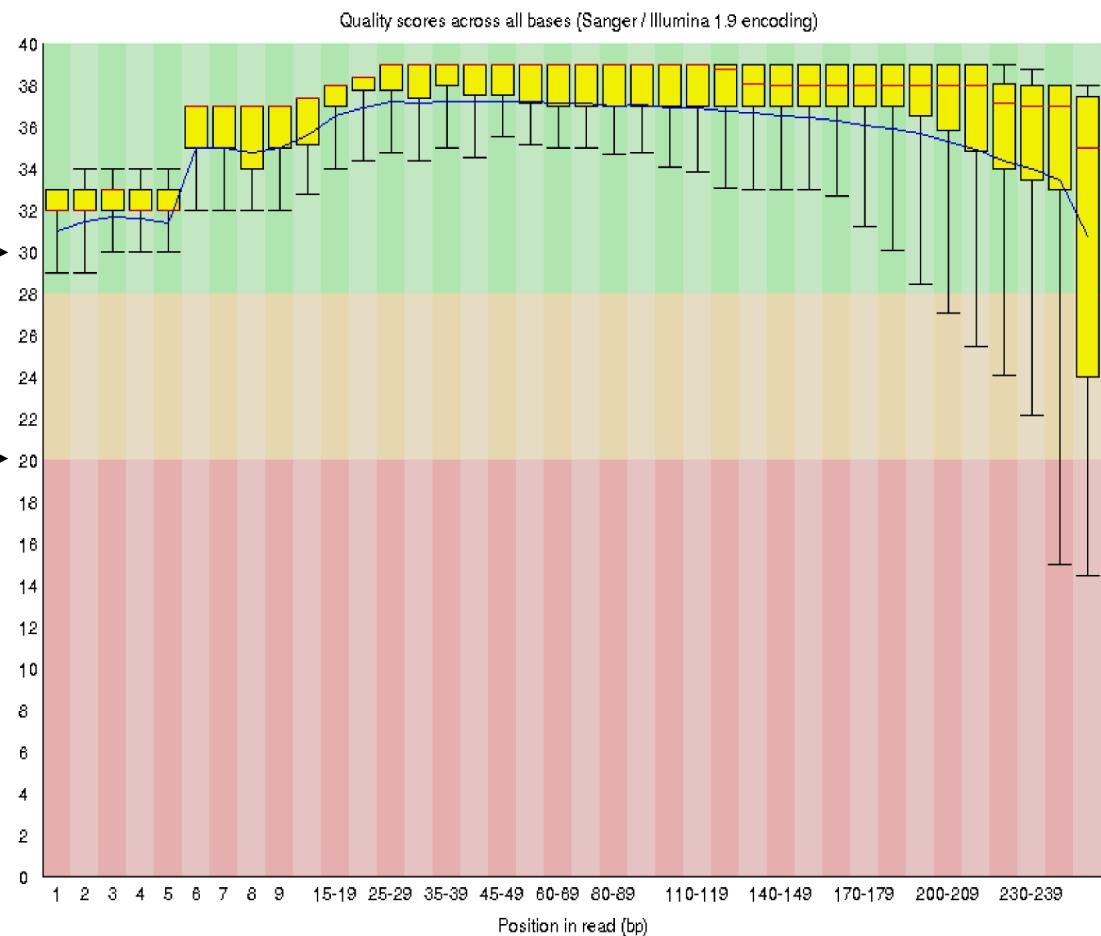
Quality check by FastQC

Raw reads → **Quality check** → Mapping → Counting → DEG analysis → Functional enrichment → Data interpretation

Accuracy

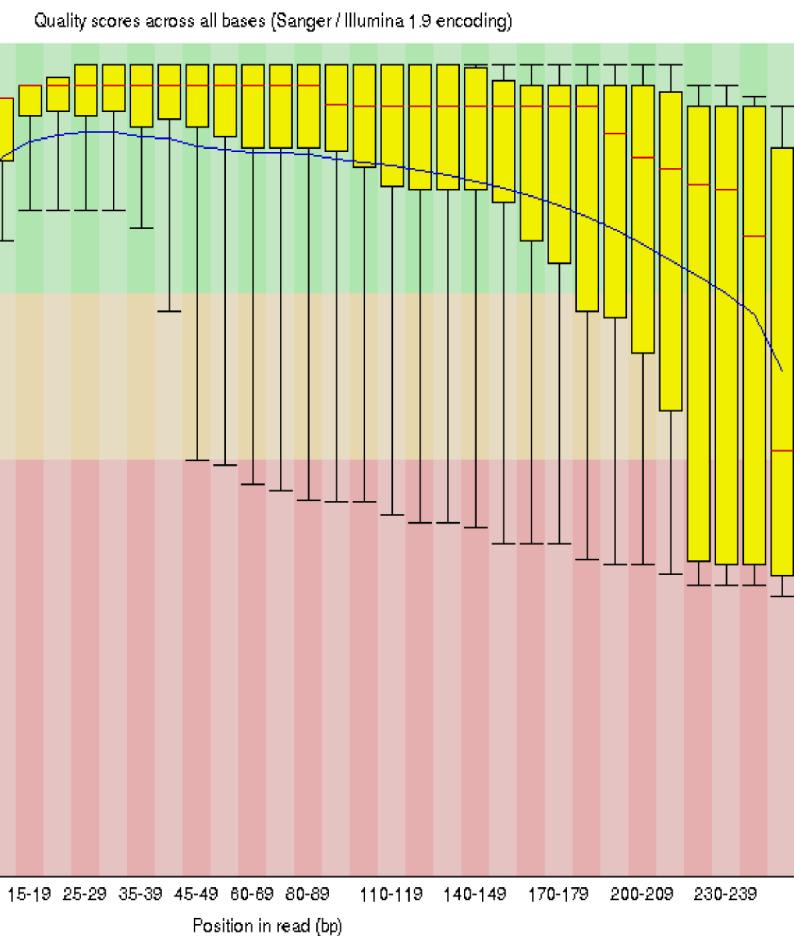
99.9% →

99% →



- 4.6 Mbp of *E. coli*
- : 1% → 46,000
- : 0.1% → 4,600

- 3,088 Mbp of Human
- : 1% → 30,880,000
- : 0.1% → 3,088,000



<https://microsizedmind.wordpress.com/2014/02/21/is-improvement-possible-quality-control-of-a-illumina-nextera-dataset-for-the-nova-genome-assembly/>

Mapping by bowtie2

Raw reads → Quality check → **Mapping** → Counting → DEG analysis → Functional enrichment → Data interpretation

10414362 reads; of these:

10414362 (100.00%) were paired; of these:

374434 (3.60%) aligned concordantly 0 times

9833851 (94.43%) aligned concordantly exactly 1 time

206077 (1.98%) aligned concordantly >1 times

374434 pairs aligned concordantly 0 times; of these:

234742 (62.69%) aligned discordantly 1 time

139692 pairs aligned 0 times concordantly or discordantly; of these:

279384 mates make up the pairs; of these:

148566 (53.18%) aligned 0 times

117701 (42.13%) aligned exactly 1 time

13117 (4.69%) aligned >1 times

99.29% overall alignment rate

Table 1. The adopted mapping methods.

Name	Version	Mapping
Bowtie	2.2.6	Unspliced read aligner
BWA 	0.7.12-r1039	Unspliced read aligner
TopHat	2.10	Spliced read aligner
STAR	2.5.3	Spliced read aligner
kallisto	0.43.1	pseudo-alignment
Salmon	0.8.2	pseudo-alignment

An example alignment result

A

Coor	12345678901234	10	20	30	40
ref	AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT				
+r001/1		TTAGATAAAGGATA*CTG			
+r002		aaaAGATAA*GGATA			
+r003		gcctaAGCTAA			
+r004			ATAGCT.....	TCAGC	
-r003			ttagctTAGGC		
-r001/2				CAGCGGCAT	

An example alignment result by SAM (Sequence Alignment/Map) format

B

QHD VN:1.5 SO:coordinate @SQ SN:ref LN:45	Header section	QUAL (read quality; * meaning such information is not available)
r001 99 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *		
r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *		
r003 0 ref 9 30 5S6M * 0 0 GCCTAAGCTAA *		
r004 0 ref 16 30 6M14N5M * 0 0 ATAGCTTCAGC *		
r003 2064 ref 29 17 6H5M * 0 0 TAGGC *		
r001 147 ref 37 30 9M = 7 -39 CAGCGGCAT *		

Alignment section

QNAME FLAG RNAME POS MAPQ CIGAR RNEXT PNEXT TLEN SEQ Optional fields in the format of TAG:TYPE:VALUE

(query template name, aka. read ID) (indicates alignment information about the read, e.g. paired, aligned, etc.) (reference sequence name, e.g. chromosome /transcript id) (1-based position) (mapping quality) (summary of alignment, e.g. insertion, deletion) (reference sequence name of the primary alignment of the NEXT read; for paired-end sequencing, NEXT read is the paired read; corresponding to the RNAME column) (Position of the primary alignment of the NEXT read in the template; corresponding to the POS column) (the number of bases covered by the reads from the same fragment. In this particular case, it's 45 - 7 + 1 = 39 as highlighted in Panel A). Sign: plus for leftmost read, and minus for rightmost read

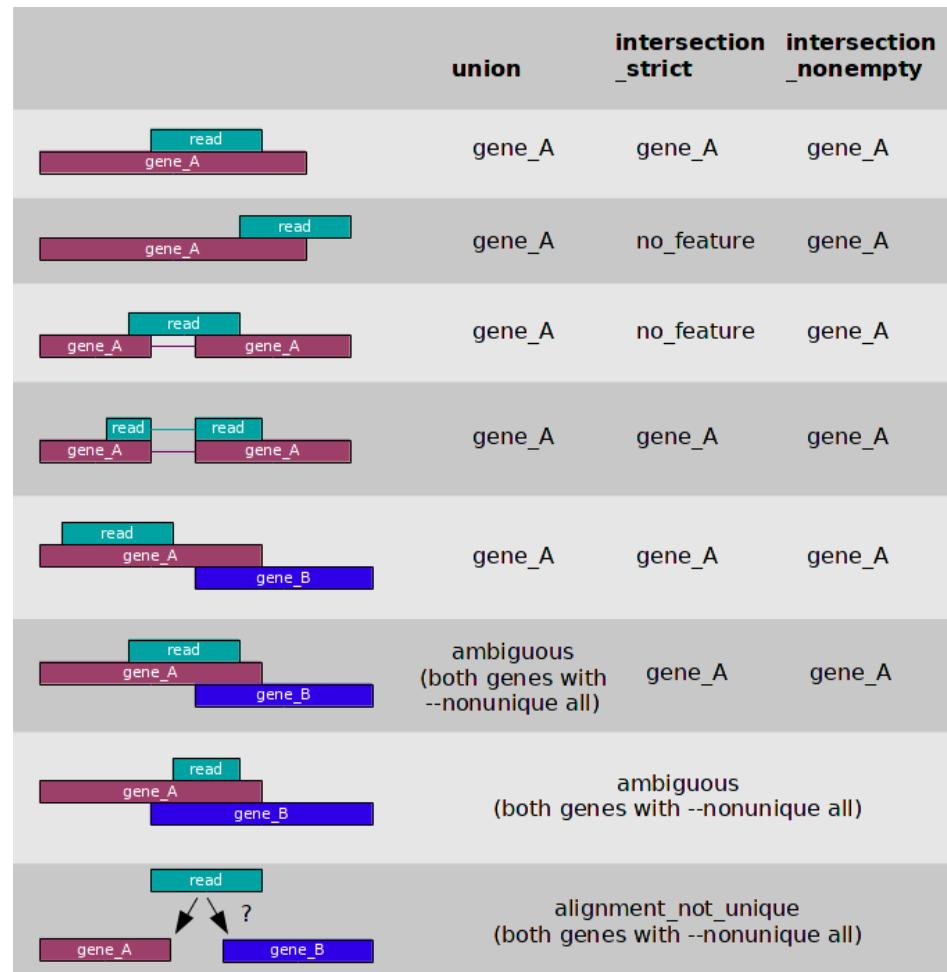
SEQ (read sequence)

*BAM: a compressed binary version of SAM

Read counting by HTseq

Raw reads → Quality check → Mapping → **Counting** → DEG analysis → Functional enrichment → Data interpretation

PSPTO_RS00025	0
PSPTO_RS00030	0
PSPTO_RS00035	0
PSPTO_RS00040	0
PSPTO_RS00045	0
PSPTO_RS00050	0
PSPTO_RS00060	0
...	
_no_feature	108
_ambiguous	162
_too_low_aQual	181
_not_aligned	126
_alignment_not_unique	0



Differentially expressed genes by DESeq2

Raw reads → Quality check → Mapping → Counting → **DEG analysis** → Functional enrichment → Data interpretation

Table 2. Adopted methods for DEGs identification.

Name	Version	Normalization
baySeq	2.4.1	Scaling factors (quantile/ TMM/ total)
DESeq	1.22.1	DESeq size factors
EBSeq	1.12.0	DESeq median normalization
edgeR	3.12.1	TMM/ Upper quartile / RLE / None (all scaling factors are set to be one)
limma+voom	3.26.9	TMM
NOIseq	2.14.1	RPKM / TMM / Upper quartile
SAMseq (samr)	2.0	Based on the read count mean over the null features of data set.
DESeq2	1.10.1	DESeq size factors
sleuth	0.29.0	DESeq size factors (with slight modifications)

Differentially expressed genes by DESeq2

Raw reads → Quality check → Mapping → Counting → **DEG analysis** → Functional enrichment → Data interpretation

Normalization method	Description	Accounted factors	Recommendations for use
CPM (counts per million)	counts scaled by total number of reads	sequencing depth	gene count comparisons between replicates of the same sample/group; NOT for within sample comparisons or DE analysis
TPM (transcripts per kilobase million)	<u>counts per length of transcript_(kb) per million reads mapped</u>	sequencing depth and gene length	gene count comparisons within a sample or between samples of the same sample group; NOT for DE analysis
RPKM/FPKM (reads/fragments per kilobase of exon per million reads/fragments mapped)	similar to TPM	sequencing depth and gene length	gene count comparisons between genes within a sample; NOT for between sample comparisons or DE analysis
DESeq2's median of ratios	<u>counts divided by sample-specific size factors</u> determined by median ratio of gene counts relative to geometric mean per gene	sequencing depth and RNA composition	gene count comparisons between samples and for DE analysis ; NOT for within sample comparisons
EdgeR's trimmed mean of M values (TMM)	uses a weighted trimmed mean of the log expression ratios between samples	sequencing depth, RNA composition, and gene length	gene count comparisons between and within samples and for DE analysis

Table 2. Adopted methods for DEGs identification.

Name	Version	Normalization
baySeq	2.4.1	Scaling factors (quantile/ TMM/ total)
DESeq	1.22.1	DESeq size factors
EBSeq	1.12.0	DESeq median normalization
edgeR	3.12.1	TMM/ Upper quartile / RLE / None (all scaling factors are set to be one)
limma+voom	3.26.9	TMM
NOIseq	2.14.1	RPKM / TMM / Upper quartile
SAMseq (samr)	2.0	Based on the read count mean over the null features of data set.
DESeq2	1.10.1	DESeq size factors
sleuth	0.29.0	DESeq size factors (with slight modifications)

RPKM-normalized counts table

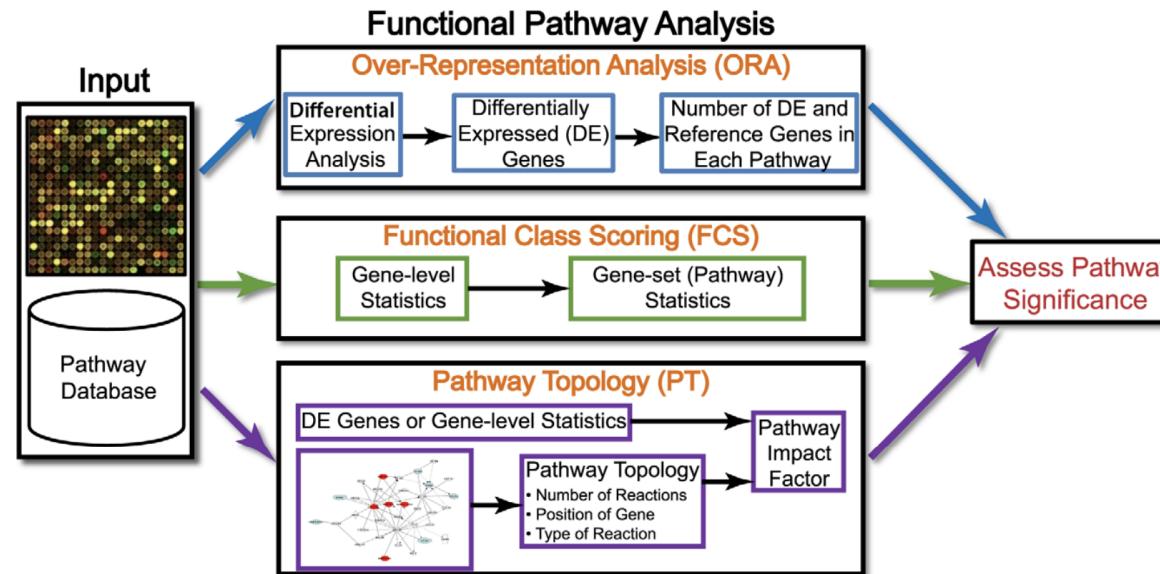
gene	sampleA	sampleB
XCR1	5.5	5.5
WASHC1	73.4	21.8
...
Total RPKM-normalized counts	1,000,000	1,500,000

DESeq2-normalized counts: Median of ratios method

gene	sample A	sample B	pseudo-reference sample	ratio of sampleA/ref	ratio of sampleB/ref	Normalized A	Normalized B
EF2A	1489	906	$\sqrt{1489 \cdot 906} = 1161.5$	$1489/1161.5 = 1.28$	$906/1161.5 = 0.78$	$1489 / 1.3 = 1145.39$	$906 / 0.77 = 1176.62$
ABCD1	22	13	$\sqrt{22 \cdot 13} = 16.9$	$22/16.9 = 1.30$	$13/16.9 = 0.77$	$22 / 1.3 = 16.92$	$13 / 0.77 = 16.88$
MEFV	793	410	$\sqrt{793 \cdot 410} = 570.2$	$793/570.2 = 1.39$	$410/570.2 = 0.72$	$793 / 1.3 = 610$	$410 / 0.77 = 532.46$
			Mean (1.28, 1.30, 1.39 ...) = 1.3			Mean (0.78, 0.77, 0.72 ...) = 0.77	

Functional enrichment

Raw reads → Quality check → Mapping → Counting → DEG analysis → **Functional enrichment** → Data interpretation



Analysis	What is required for input	What output looks like	Pros	Cons
<u>ORA (Over-representation Analysis)</u> <i>e.g., Gene ontology (GO), eggNOG-mapper, Pathview</i>	A list of gene IDs (no stats needed)	A per-pathway hypergeometric test result	Pros <ul style="list-style-type: none">- Simple- Inexpensive computationally to calculate p-values	Cons <ul style="list-style-type: none">- Requires arbitrary thresholds and ignores any statistics associated with a gene- Assumes independence of genes and pathways
<u>GSEA (Gene Set Enrichment Analysis)</u>	A list of genes IDs with gene-level summary statistics	A per-pathway enrichment score	<ul style="list-style-type: none">- Includes all genes (no arbitrary threshold!)- Attempts to measure coordination of genes	<ul style="list-style-type: none">- Permutations can be expensive- Does not account for pathway overlap- Two-group comparisons not always appropriate/feasible
<u>GSVA (Gene Set Variation Analysis)</u>	A gene expression matrix (like what you get from refine.bio directly)	Pathway-level scores on a per-sample basis	<ul style="list-style-type: none">- Does not require two groups to compare upfront- Normally distributed scores	<ul style="list-style-type: none">- Scores are not a good fit for gene sets that contain genes that go up AND down- Method doesn't assign statistical significance itself- Recommended sample size n > 10

https://alexlemonade.github.io/refinebio-examples/03-rnaseq/pathway-analysis_rnaseq_01_ora.html

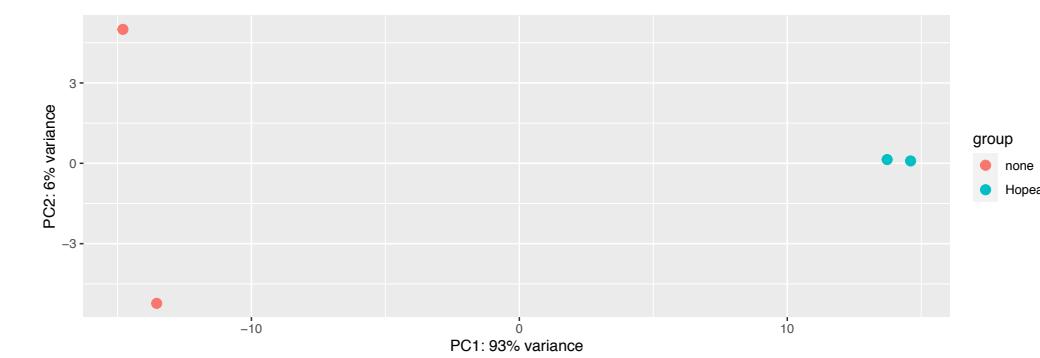
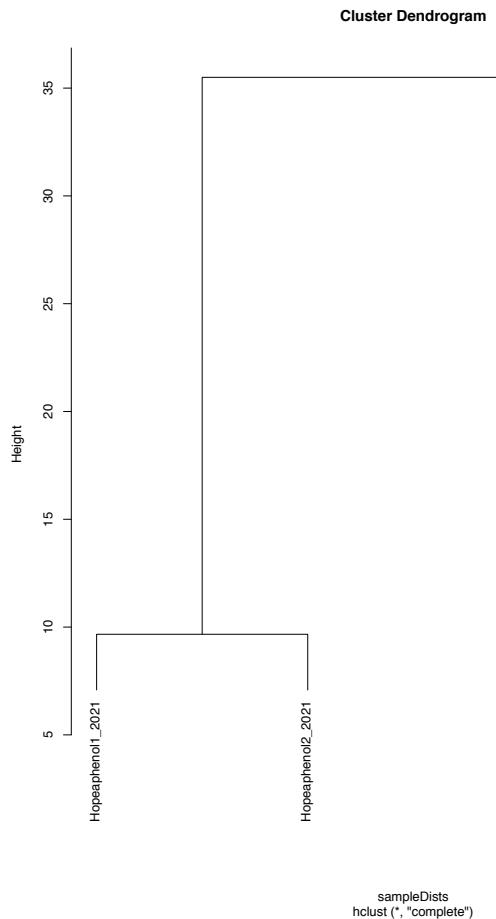
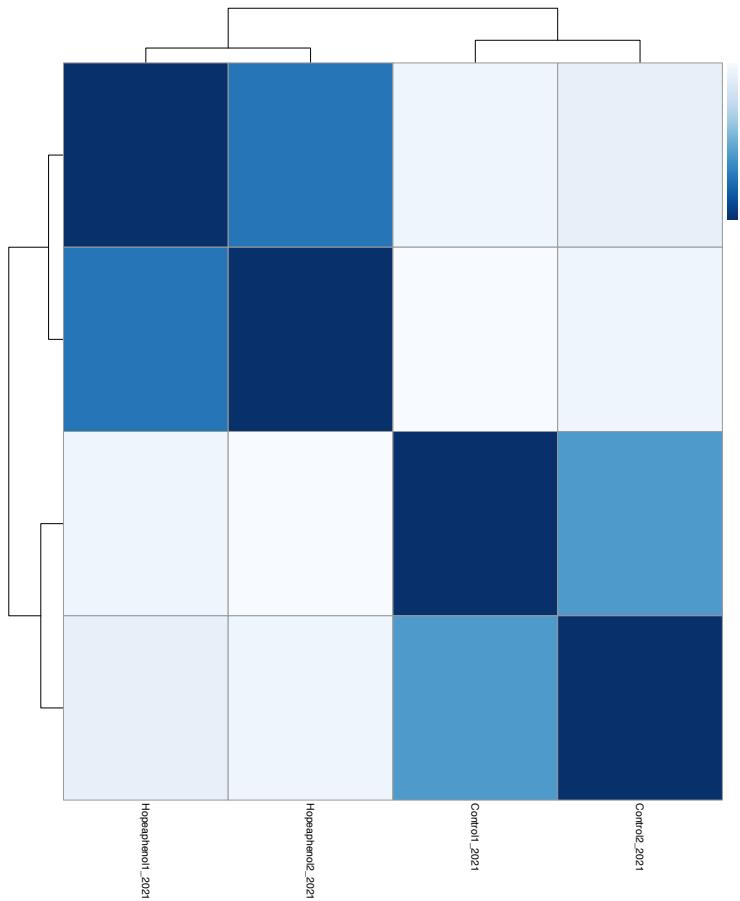
Costa-Silva et al., 2017 PlosOne

<https://doi.org/10.1371/journal.pcbi.1002375>

Transcriptomics: Effect of hopeaphenol to *Pectobacterium atrosepticum* SCRI1043

Data analysis

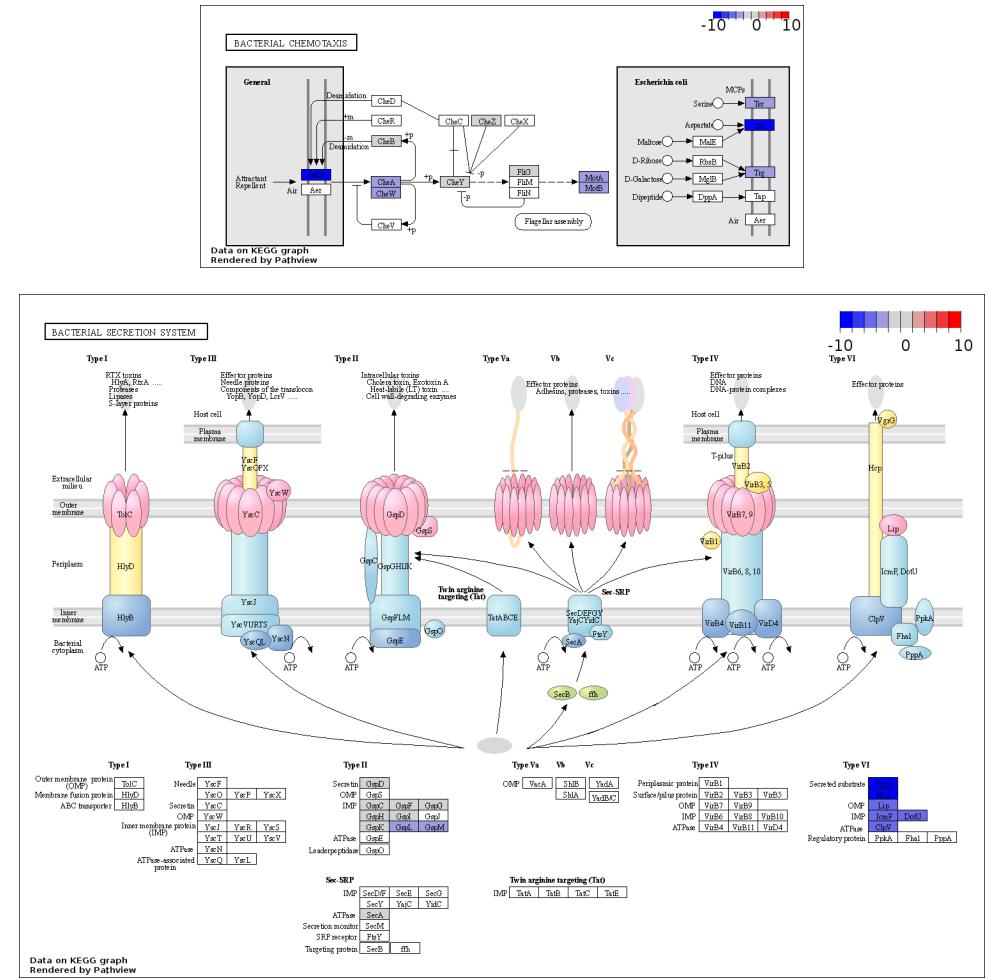
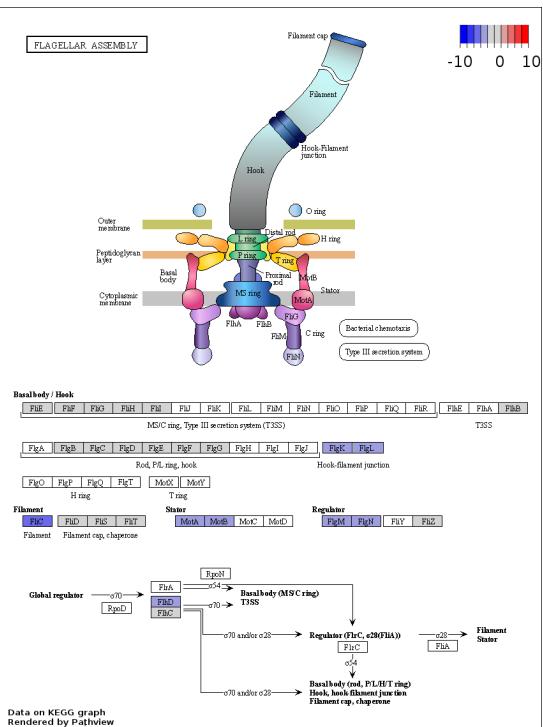
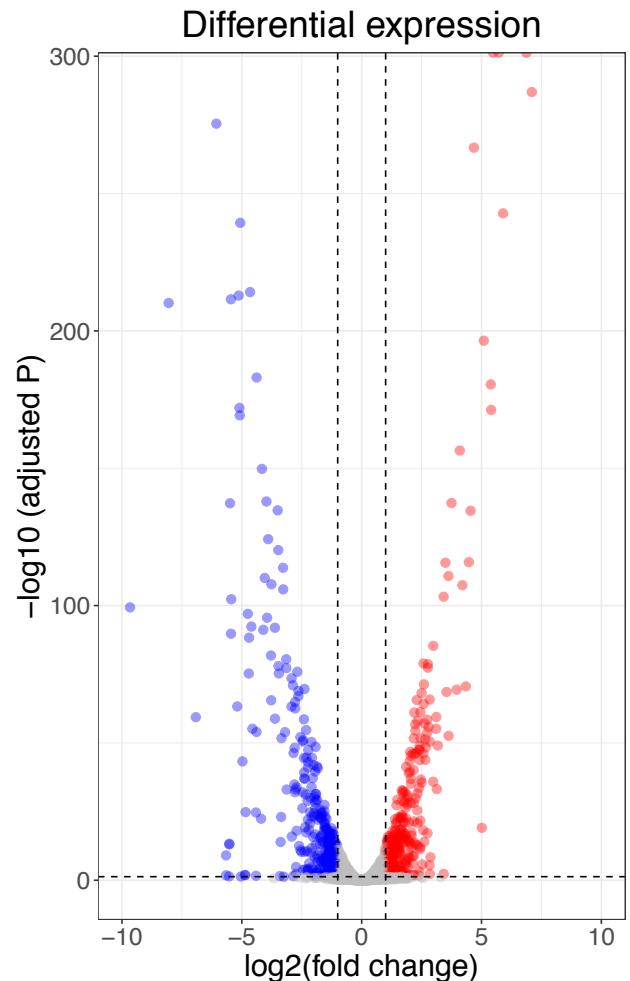
- Alignment: Bowtie2
- Read counting: HTSeq
- Statistics: DESeq2



Transcriptomics: Effect of hopeaphenol to *Pectobacterium atrosepticum* SCRI1043

Hopea vs non-treated

689 DEGs (336 up, 353 down)



Hand-on RNA-seq data analysis

GitHub repository

https://github.com/sungminhwang-duke/RNAseq_hand_on_work

The screenshot shows the GitHub repository page for 'sungminhwang-duke / RNAseq_hand_on_work'. The repository is public and has 1 branch and 0 tags. The 'Code' tab is selected. The repository's history shows several commits from 'sungminhwang-duke' updating the README.md file and adding files via upload. The 'About' section indicates no description, website, or topics provided. The 'Releases' section shows no releases published, with a link to 'Create a new release'. The 'Packages' section shows no packages published, with a link to 'Publish your first package'. The 'Languages' section is currently empty.

Code

main · 1 branch · 0 tags

Go to file Add file Code

sungminhwang-duke Update README.md 9b82159 16 hours ago 18 commits

DEGs Add files via upload 17 hours ago

Raw_materials Add files via upload 17 hours ago

Setup_env Add files via upload 17 hours ago

README.md Update README.md 16 hours ago

README.md

RNAseq_hand_on_work

1. Go to the folder "Seup_env" and download a file to make a virtual environment for a data process.
2. Execute this command using Terminal (Mac only): `conda env create -f condaTESTmac.yml`
-Keep in mind that the path should be identical to the Terminal and the downloaded file, "condaTESTmac.yml".
3. Download raw data (toy data) and materials for the data process.
4. For the DEG analysis, go to "DEGs" folder and download the exercise dataset (toy data) and DESeq2 command