# Biomarkers for disease identification/outcome

INDRAPRASTHA INSTITUTE *of* INFORMATION TECHNOLOGY **DELHI**

Dr. Jaspreet Kaur Dhanjal

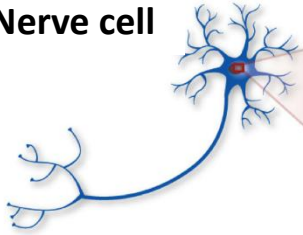Assistant Professor, Department of Computational Biology
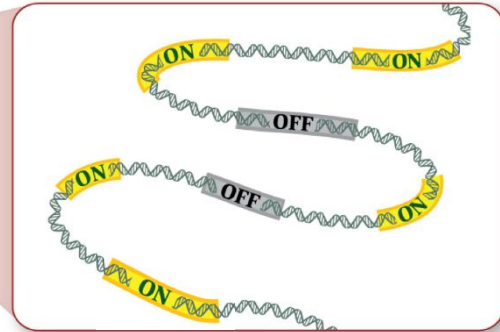
Email ID: jaspreet@iiitd.ac.in

*October 14, 2025*

1

# Why Transcriptome?
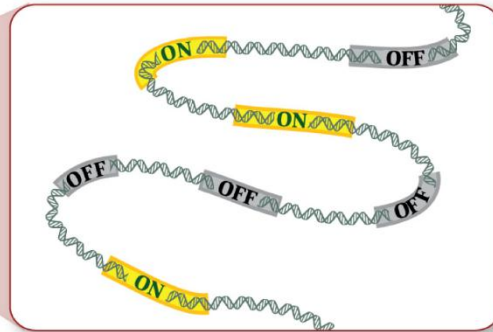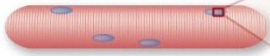
# Differential gene expression analysis



which genes are expressed at different levels and reasonable for the disease ?

Normal cell    VS    Tumor cell

Gene A    Gene C
Gene B
Gene D

Gene A    Gene C
Gene B
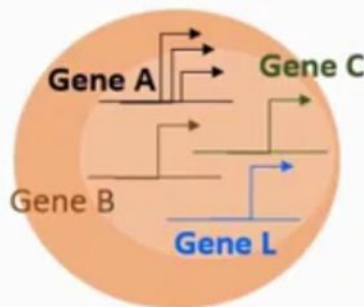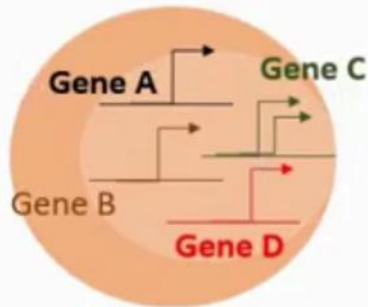Gene L

What are the differentially expressed genes?

Gene A is up regulated
Gene c is down regulated
Gene D is turned off
Gene L is turned on

# Why sequence RNA (Versus DNA)?

1. *Functional studies*

Genome may be constant but an experimental condition has a profound effect on the gene expression (differential expression)

Eg. Drug *vs.* untreated cells

Eg. Wild type *vs.* knock out mice cells

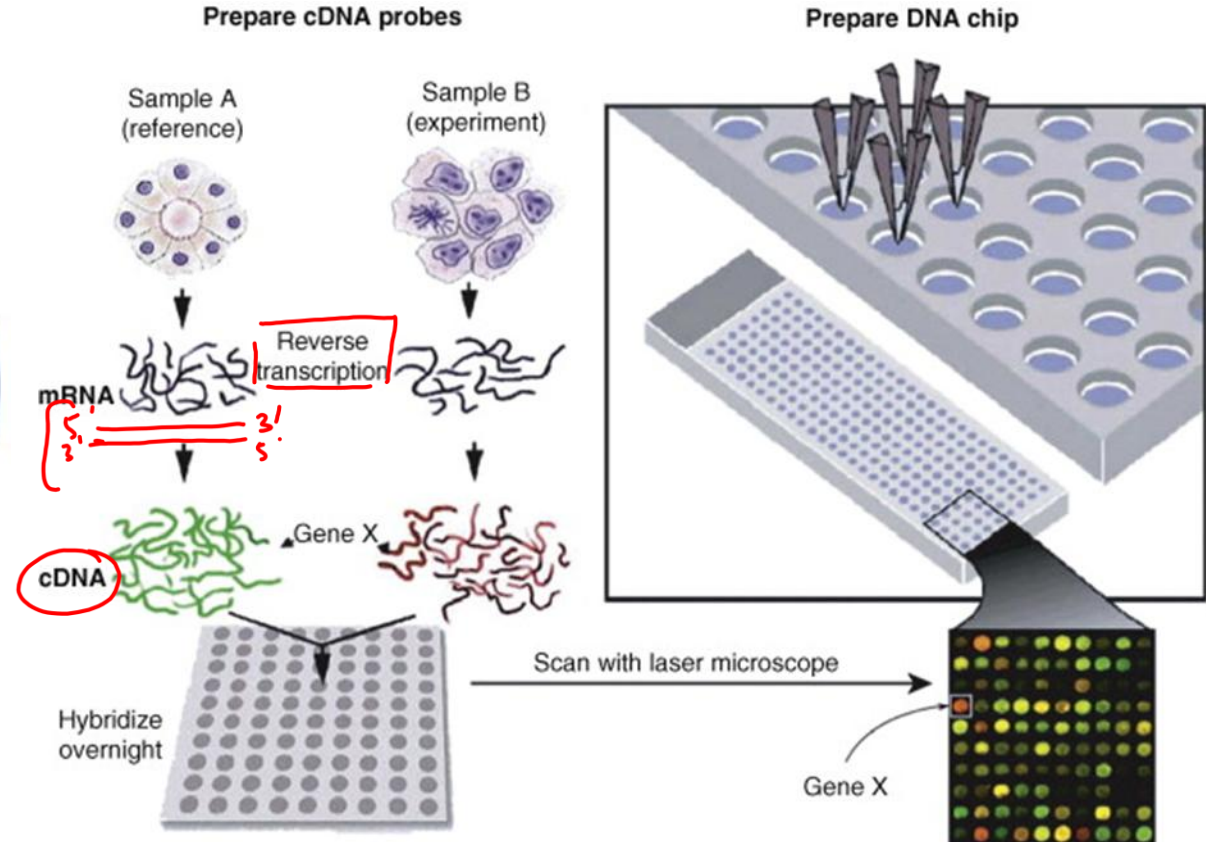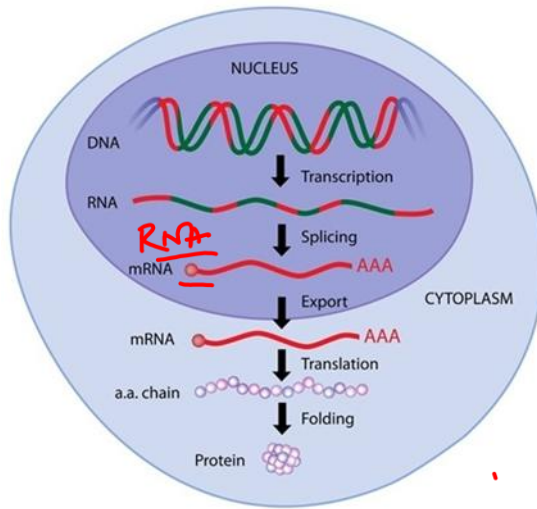2. *Predicting transcript sequence from genome sequence is difficult*

3. *Some molecular features can only be observed at the RNA level*
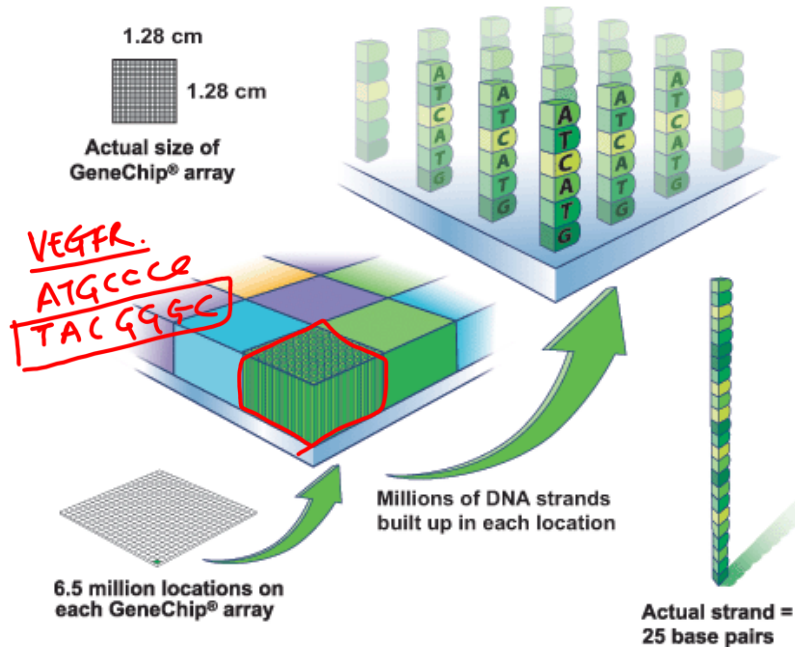
Alternative isoforms, fusion transcripts, RNA editing

4. *Understand allele specific expression*

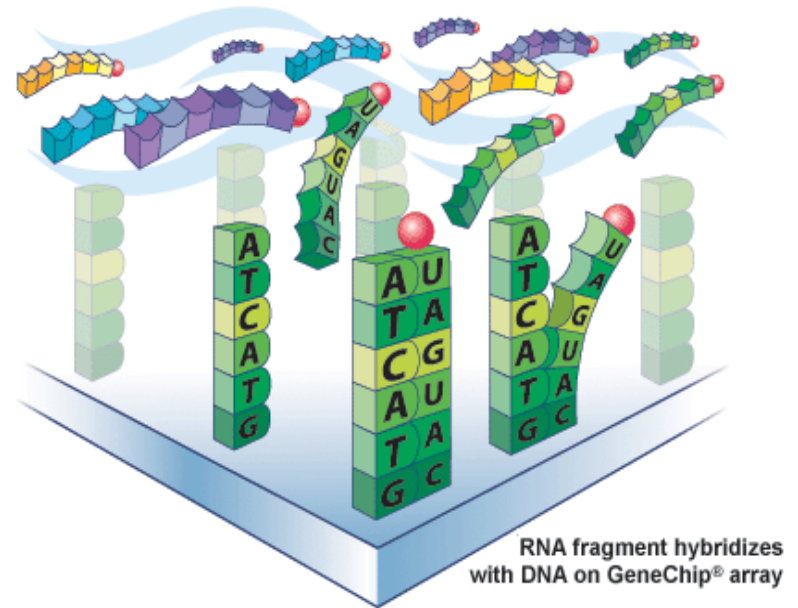# Gene expression-based biomarker identification

**DNA Microarray**
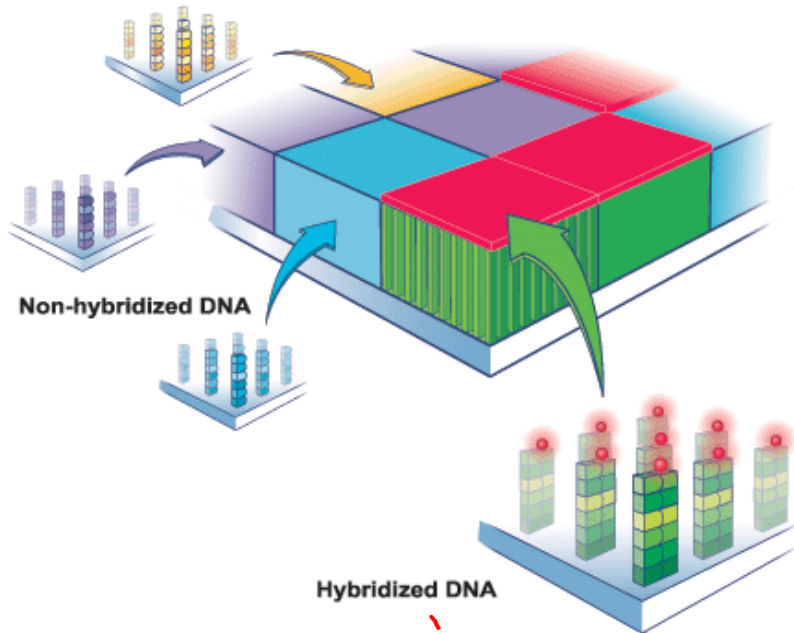
# Microarray analysis



**1. Micro Array Features**

**2. Hybridization (Pairing)**

# Microarray analysis

Shining a laser light at GeneChip® array causes tagged DNA fragments that hybridized to glow
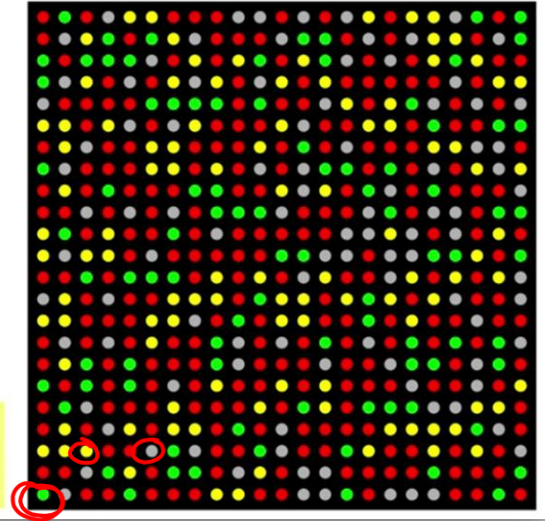
Non-hybridized DNA

Hybridized DNA

**3. Detection**

cDNAs from tissue 1 were labeled red

cDNAs from tissue 2 were labeled green

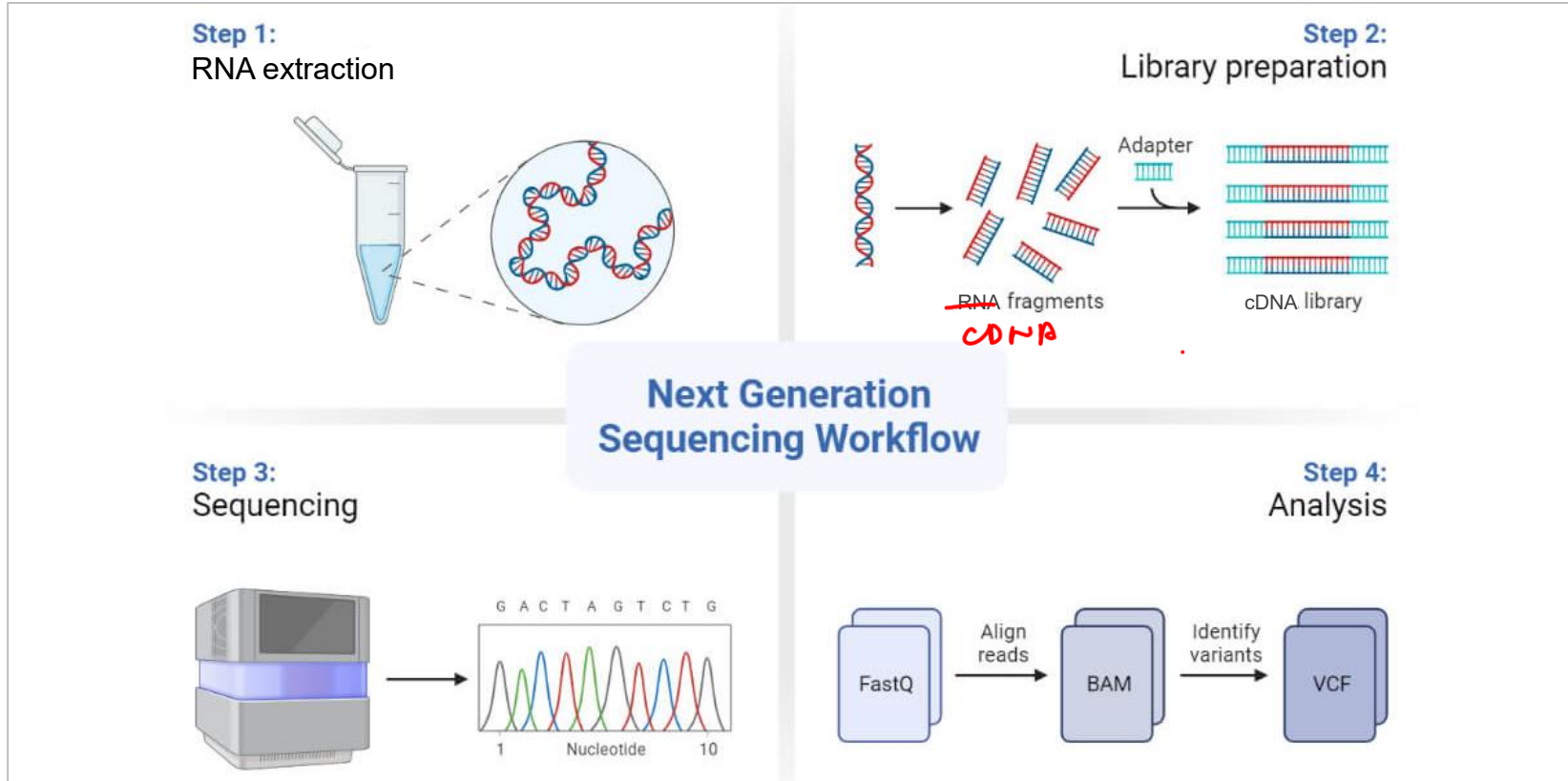red spot means gene is expressed in tissue 1

green spot means gene is expressed in tissue 2

yellow spot means both cDNAs bind and gene is expressed in both tissues

**Limitations:**
1. Relies on existing knowledge about genome sequence.
2. Technical problems like high background levels owing to cross-hybridization
3. Comparison of expression across different samples/experiments is often complicated.
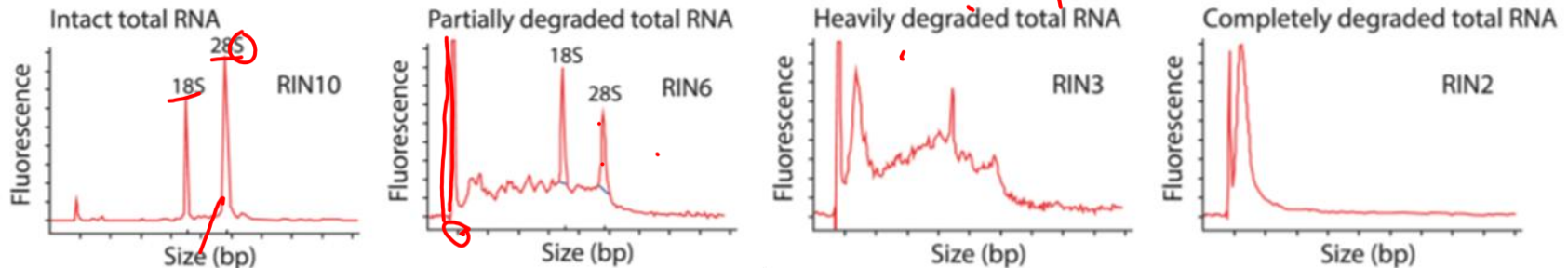
# Next generation RNA-sequencing

## Types of RNA

poly-adenylated (coding) RNAs, "genes"

short non-coding RNAs (ncRNA), "microRNA"  } Total RNA

long non-coding RNAs

ribosomal RNA

*rRNA*

Most of the RNA in the cell { polyA capture, ribominus } "Enrichment"

*RNA size*

(a) Gel electropherogram

Intact total RNA | Partially degraded total RNA | Heavily degraded total RNA | Intact mRNA

bp 12,000, 5,000, 2,000, 1,650, 1,000, 850, 650, 500, 400, 300, 200, 100

(b) Capillary electropherogram

Intact total RNA — Fluorescence / Size (bp) — 18S, 28S, RIN10

Partially degraded total RNA — Fluorescence / Size (bp) — 18S, 28S, RIN6

Heavily degraded total RNA — Fluorescence / Size (bp) — RIN3

Completely degraded total RNA — Fluorescence / Size (bp) — RIN2

# Library Preparation

# RNA-seq experiment workflow

**Library preparation**

# Illumina sequencing platforms



MiniSeq   MiSeq   NextSeq   HiSeq   HiSeq X   NovaSeq

Benchtop

Production-Scale

HiSeq 2500 (2 lane)   HiSeq 3000/4000   NextSeq 500

***Other sequencing platforms*:** Pacific Bioscience,  Oxford Nanopore, 10X Genomics

# Sequencing by synthesis



DNA
(0.1-5.0 µg)

**Library Preparation**

*Single molecule array*

**Cluster Growth**

3' 5'

5'

**Sequencing**

1  2  3  4  5  6  7  8  9

**Image Acquisition**

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

**Base Calling**

**Illumina sequencing video URL: https://www.youtube.com/watch?v=womKfikWlxM**

13

# Single- and Paired-end sequencing

# Multiplexing

# Microarray *vs.* RNA-seq

| Microarray | RNA-seq |
|---|---|
| • Limited probe-set based on prior knowledge of the transcriptome | • Comprehensive overview of the transcriptome |
| • Higher throughput | • Best dynamic range |
| • Analysis is more user-friendly than RNA-seq currently | • Gene fusion, isoform, SNPs detection |

# RNA-seq experiment workflow

**Sample preparation**



1. *Biological replicates* : Include multiple sampling within the population
2. *Technical replicates* : Include multiple preparation and re-sequencing of the same sample

Biological replicates generally increase statistical power more than technical replicates
- Biological variability is generally greater than technical variability
- Biological replicates contain both biological and technical variability

# Sources for RNA-seq datasets

# Sources for RNA-seq datasets

# Sources for RNA-seq datasets

# Sources for RNA-seq datasets

# Sources for RNA-seq datasets

# Workflow of differential gene expression analysis

# Gene expression-based biomarker identification

**Problems in sequencing**

1. Low confidence bases, Ns
2. Specific sequence bias, GC bias
3. Adaptors
4. Sequence contamination



NOTE: for paired-end runs, there is a second file
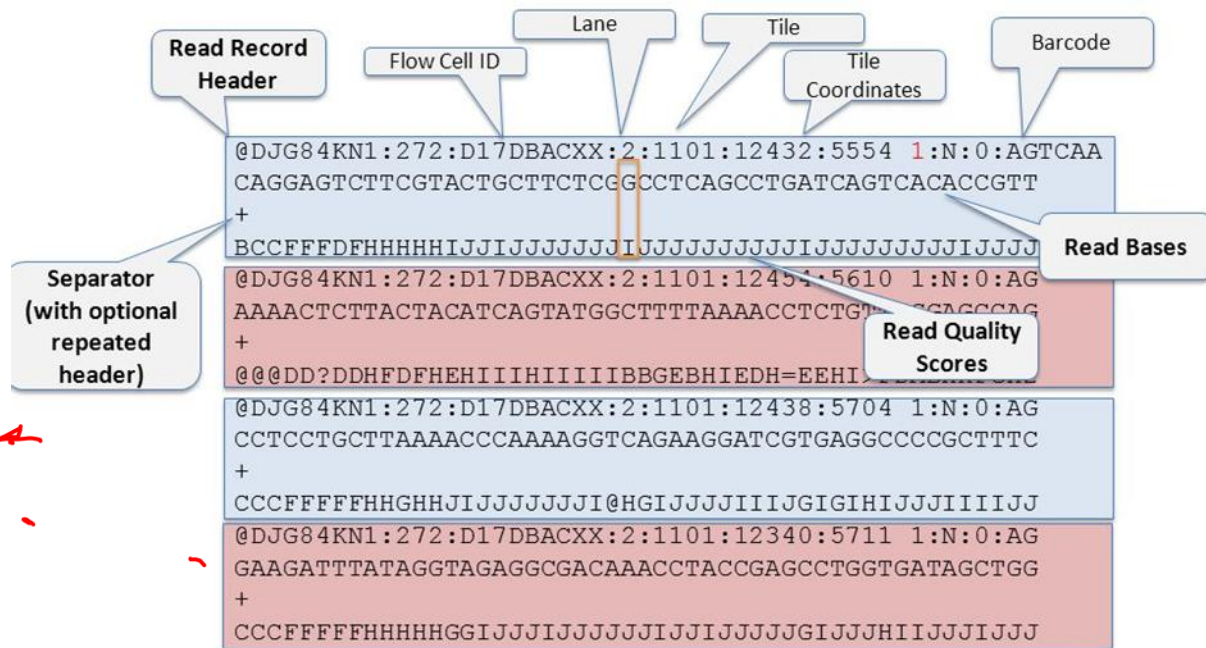with one-to-one corresponding headers and reads.

# Gene expression-based biomarker identification

## Problems in sequencing

1. Low confidence bases, Ns
2. Specific sequence bias, GC bias
3. Adaptors
4. Sequence contamination



Phred quality scores are logarithmically linked to error probabilities

| Phred Quality Score | Probability of incorrect base call | Base call accuracy |
|---|---|---|
| 10 | 1 in 10 | 90% |
| 20 | 1 in 100 | 99% |
| 30 | 1 in 1000 | 99.9% |
| 40 | 1 in 10,000 | 99.99% |
| 50 | 1 in 100,000 | 99.999% |
| 60 | 1 in 1,000,000 | 99.9999% |