

BMEG 3105 Fall 2024

Data analytics for personalized genomics and precision medicine

Lecturer: Yu LI (李煜) from CSE

Lecture 16: Cancer genomics overview & genomics analysis (30/10/2024)

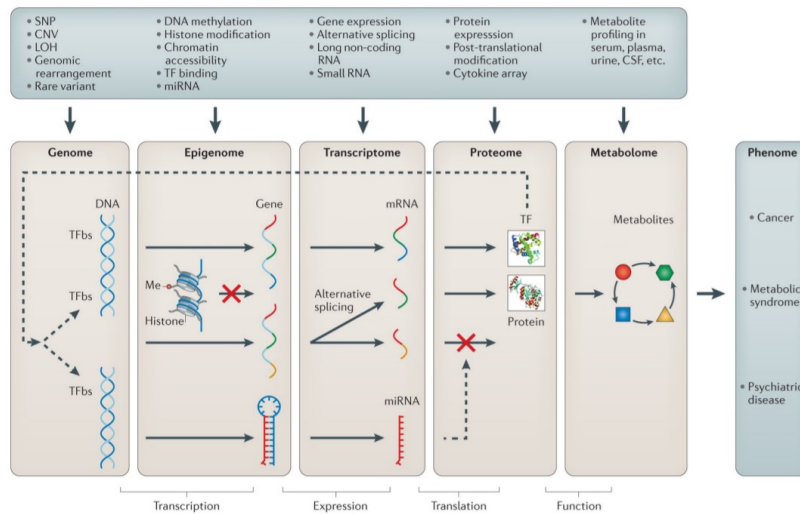
Scriber: Tom Yuet Yi (1155160583)

Cancer

- Disease with uncontrollable cell growth
- Symptom of spreading to other body parts
- Often identified as a genomic disease

➔ Study through genomics and multi-omics.

Multi-omics level: genome → epigenome → transcriptome → proteome → metabolome



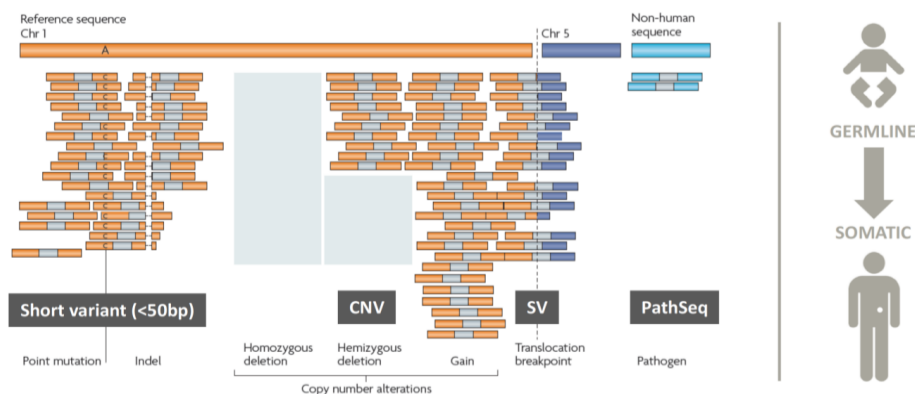
Data Analytics for Cancer Genomics

Genome: variant calling, genome-wide association studies (GWAS)

- Variant calling

Genomic variant: difference of a genome with relation to a reference.

Different types of genomic variants



➔ Contribute to: disease risk, response to treatment, variants that contribute to disease

➔ Cancer have genetic variants at multiple levels.

Genetic variant or Errors?

Genetic variants: the actual change // Errors: artifacts that introduce to the analysis.

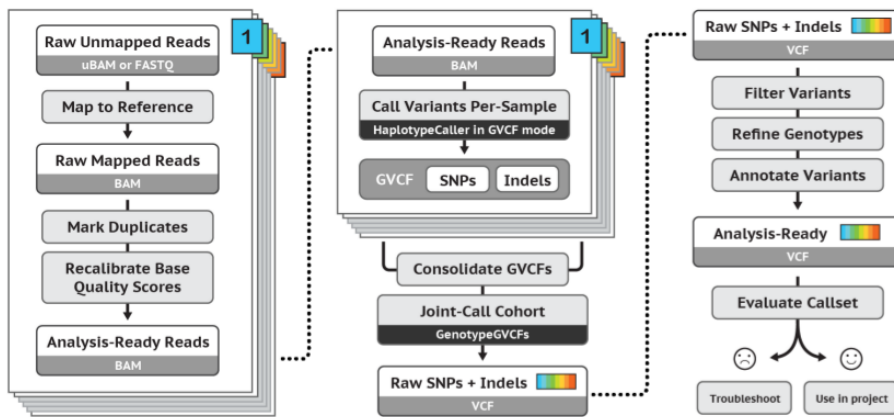
Errors can be found at:

- PCR artifacts
- Sequencing
- Alignment
- Variant calling
- Genotyping

Steps of discovering the genetic variants:

Enormous pile of short reads in terms of uBAM or FASTQ (Library preparation, Sequencing)

→ Reads mapped and cleaned up (Data pre-processing steps: 1) mapping, 2) marking duplicates, 3) base recalibration) → obtain the list of variants



1) Map the reads produced to reference:
 input as uBAM or FASTQ, output as SAM or BAM
 Alignment structure summarised by CIGAR

2) Mark duplicates to mitigate duplication artifacts (duplicates are non-independent measurements of a sequence fragment and must be removed for alleles correctly)
 **Where does duplicates from?? → Library (during PCR) or optical duplicates (during sequencing)

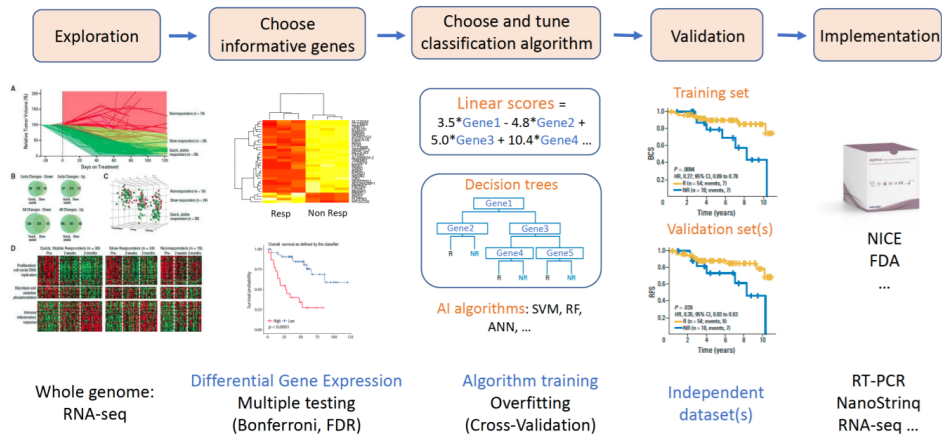
- Genome-Wide Association Studies (GWAS)
 aim to identify associations between genetic variants and traits or diseases.
 The sheer number of SNPs (e.g., 3.5 million) complicates the detection of significant associations, solution? → Bonferroni correction: Adjusted p-values are calculated to account for multiple testing, reducing the chance of false positives.

❖ Adjusted p-value = p-value/number of tests

❖ Suppose we have 1 million SNPs to test

- Adjusted p-value = $\frac{0.05}{1,000,000}$
- Adjusted p-value = $5 * 10^{-8}$

1. -Seq Data Analysis



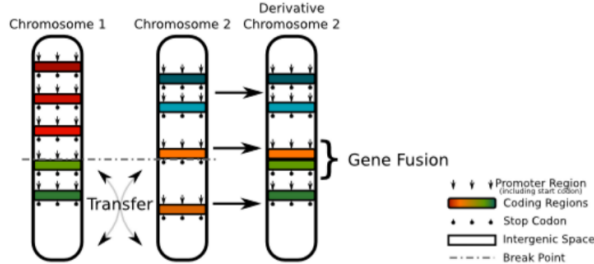
Gene Fusion – a specific kind of structural variant related to cancer

the first fusion gene was described in cancer cells

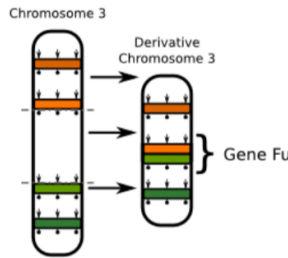
Novel gene formed by fusion of 2 distinct wild type genes

in cancer: produced by somatic genome rearrangement

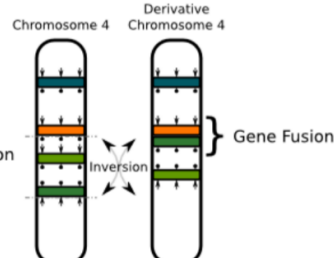
A. Chromosomal Translocation



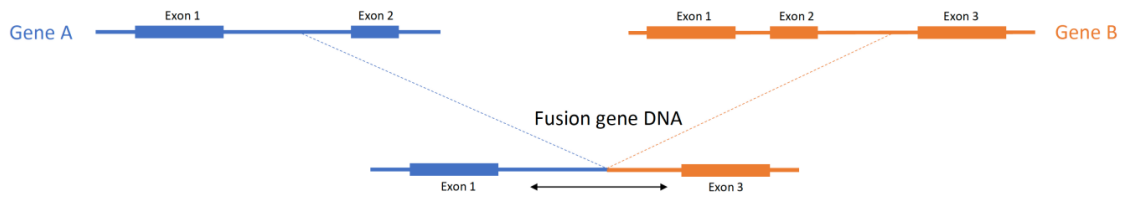
B. Interstitial Deletion



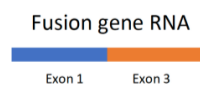
C. Chromosomal Inversion



RNA-seq for gene fusion detection

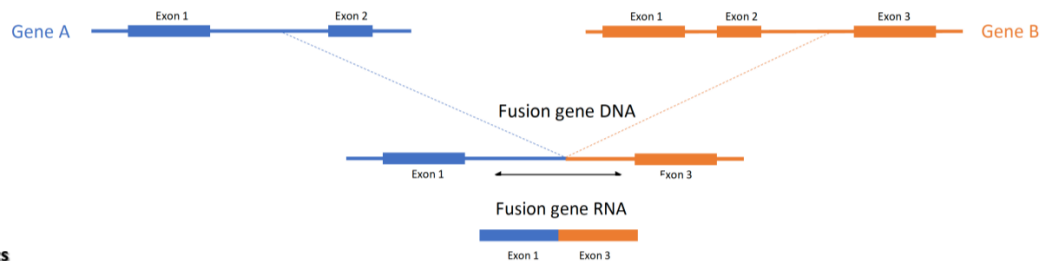
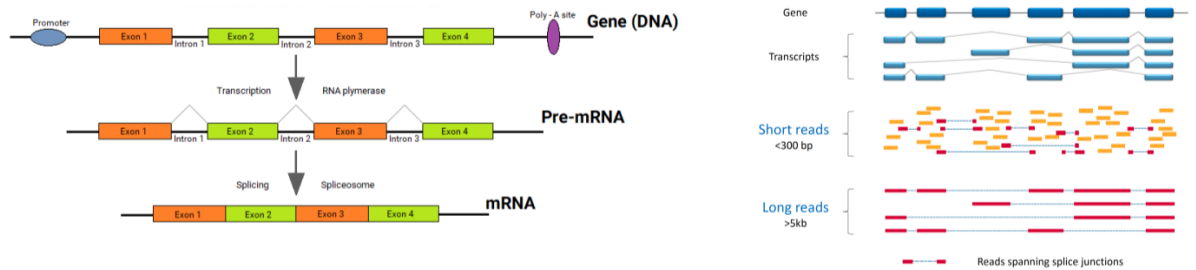


Break-points are in **introns**
 We need **whole genome sequencing**
 Whole exome sequencing is not enough



Detecting fusion in **RNA-seq** requires much less sequencing than WGS, especially with long reads

Why can it be detected by RNA-seq?



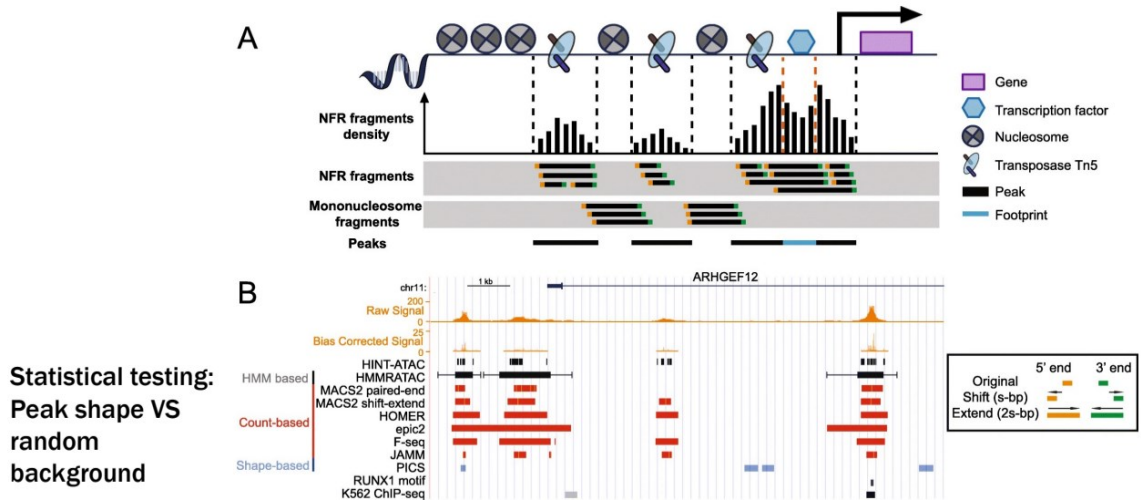
inccer genomics

2. Epigenomics

Peak Calling:

- Peak calling is a key step in the epigenomics data analysis pipeline.
- It aims to identify regions of the genome with significantly higher signal compared to the background.
- This involves statistical testing to determine which peaks are genuine vs. random background.
- statistical testing is used to compares the peak shape to the random background signal.

- This allows distinguishing true peaks from noise.

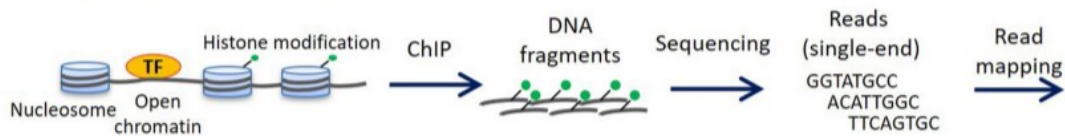


Peak Calling Output:

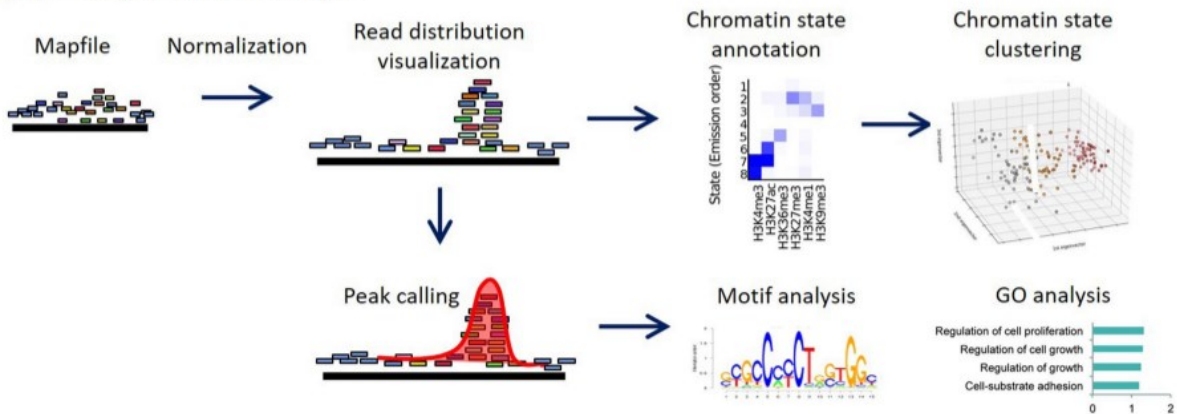
- The output of peak calling is typically in Browser Extensible Data (BED) format.
- This includes the chromosome, start, end, and label/annotation for each identified peak.

The Epigenomics Data Analysis Pipeline:

(A) Sample preparation and sequencing



(B) Computational analysis



The Detailed Epigenomics Pipeline:

