



# BIO214 Lecture 11

**Bioinformatics-II**

***Genome Assembly & variant analysis***

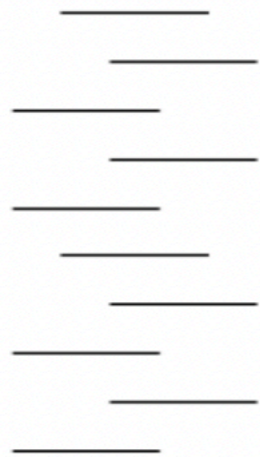
Zhen Wei; 2023-Feb-17

# Outline

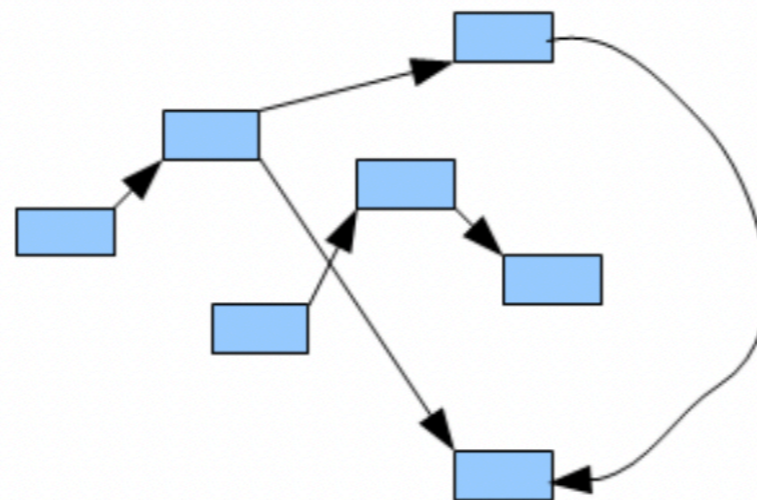
- Efficient assembly of short reads with De Bruijn graph
- Detection of genetic variants
- GWAS
- eQTL

# An efficient algorithm for short reads based genome assembly

NGS library



de Bruijn Graph



Genome



De Bruijn graph-based genome assembly algorithm:

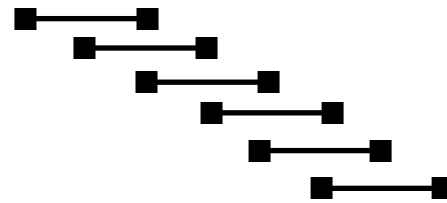
- Step 1: Short reads broken into small pieces (k-mers) and de Bruijn graph constructed.
- Step 2: Genome derived from de Bruijn graph by finding the longest possible path (Eulerian walks).

# De Bruijn Graph

Genome: AAABBBBA

# De Bruijn Graph

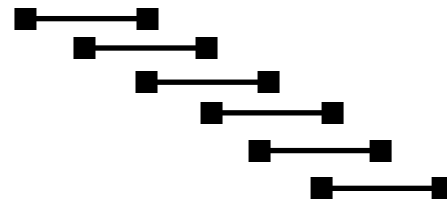
Genome: AAABBBBA



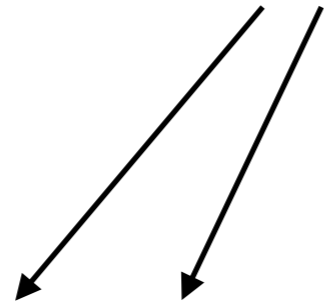
3-mers: AAA AAB ABB BBB BBB BBA

# De Bruijn Graph

Genome: AAABBBBA



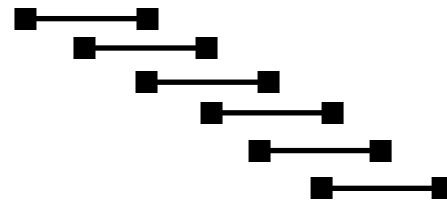
3-mers: AAA AAB ABB BBB BBB BBA



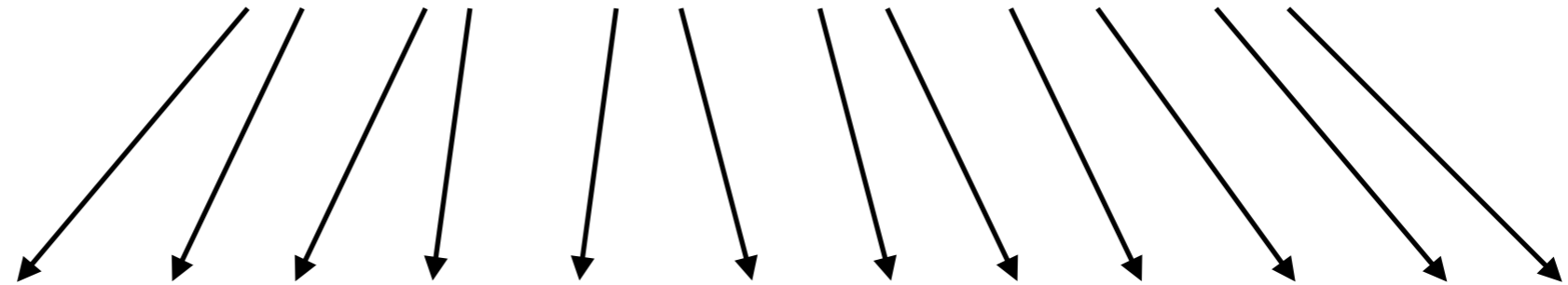
L-R 3-1mers: AA, AA

# De Bruijn Graph

Genome: AAABBBBA



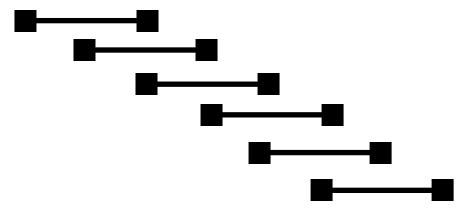
3-mers: AAA AAB ABB BBB BBB BBA



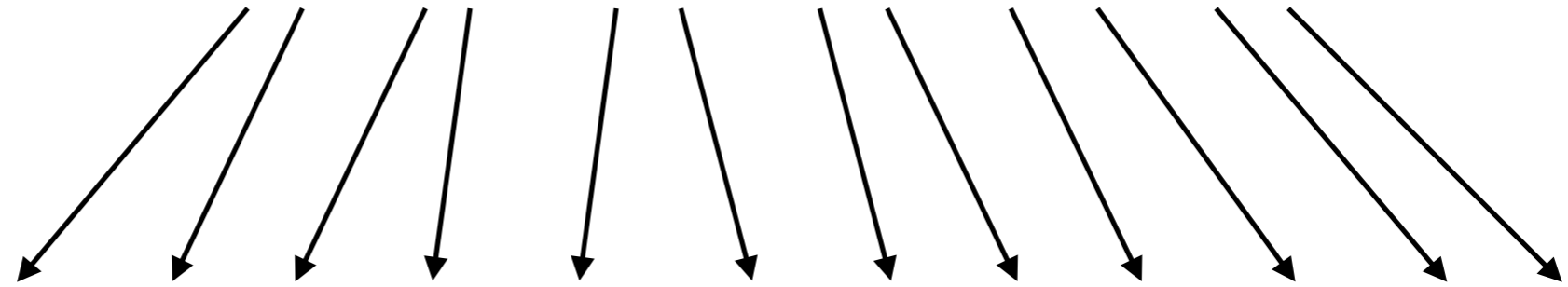
L-R 3-1mers: AA, AA AA, AB AB, BB BB, BB BB, BB BB, BA

# De Bruijn Graph

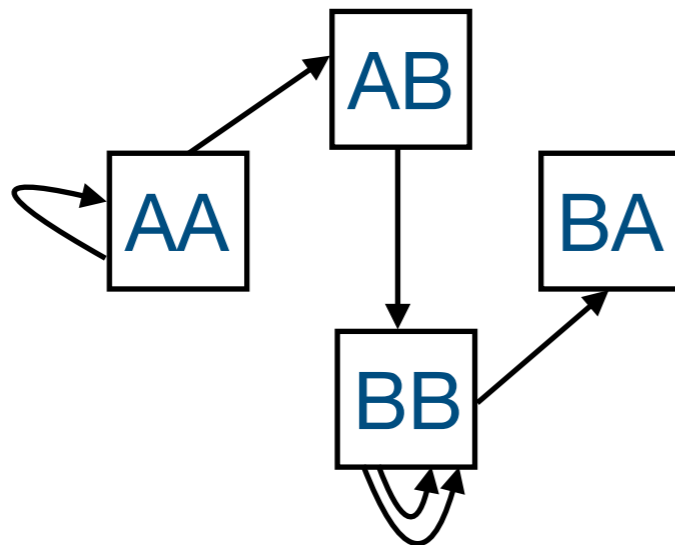
Genome: AAABBBBA



3-mers: AAA AAB ABB BBB BBB BBA



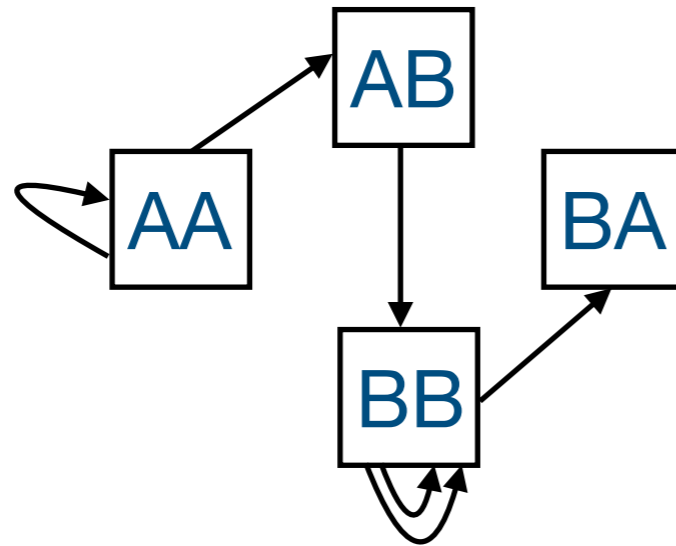
L-R 3-1mers: AA, AA AA, AB AB, BB BB, BB BB, BB BB, BA



- One edge per k-mer
- One node per distinct k-1 mer

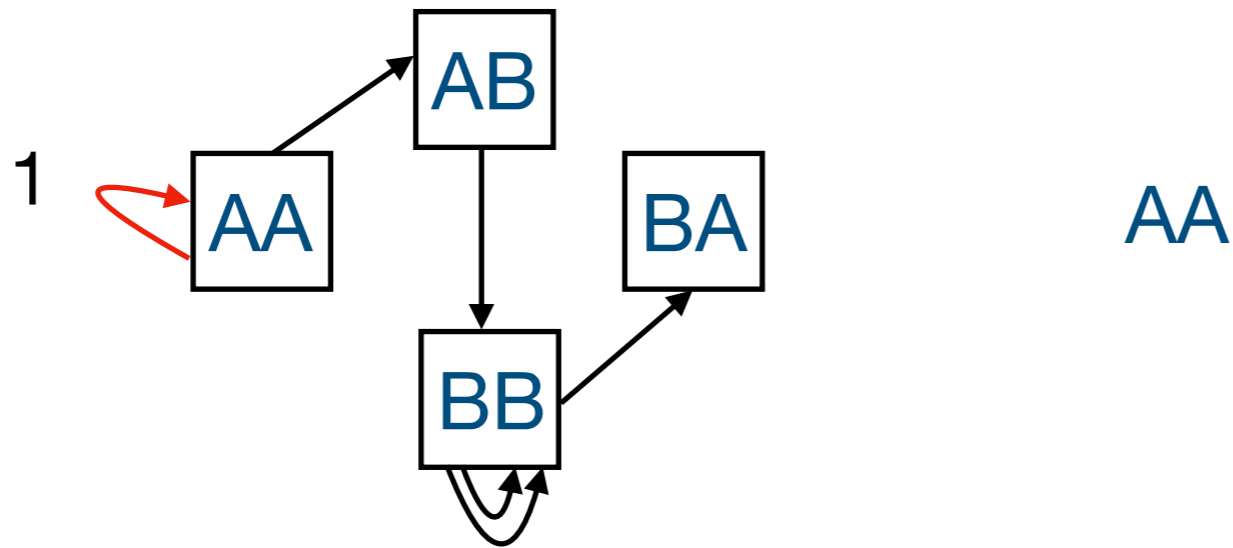


# Eulerian walk



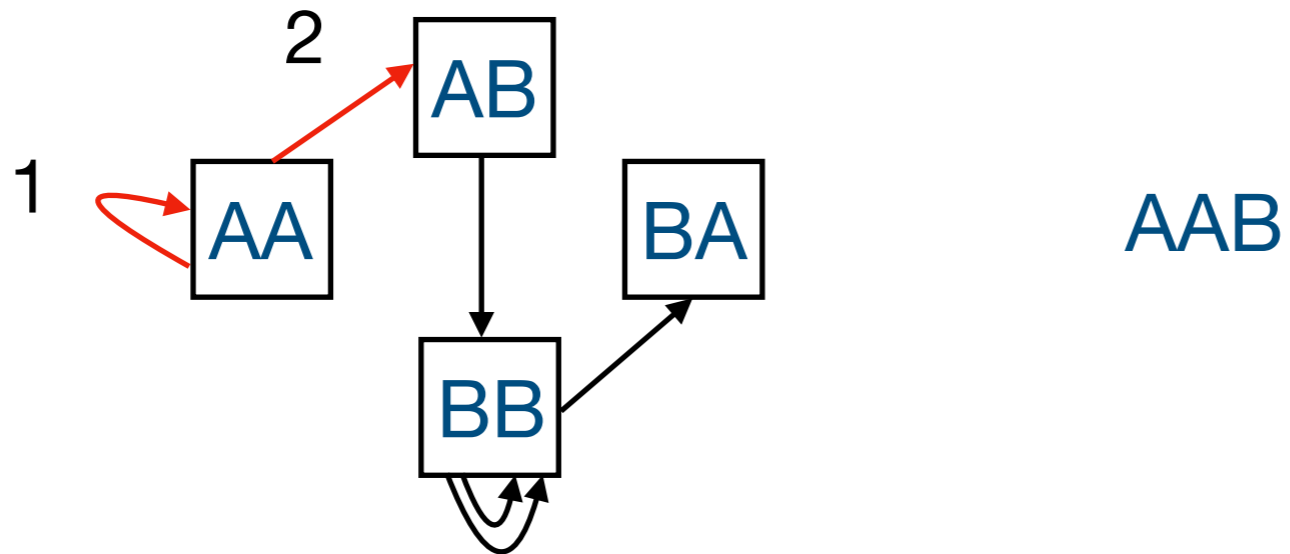
- Walk crossing each edge exactly once gives a reconstruction of the genome.

# Eulerian walk



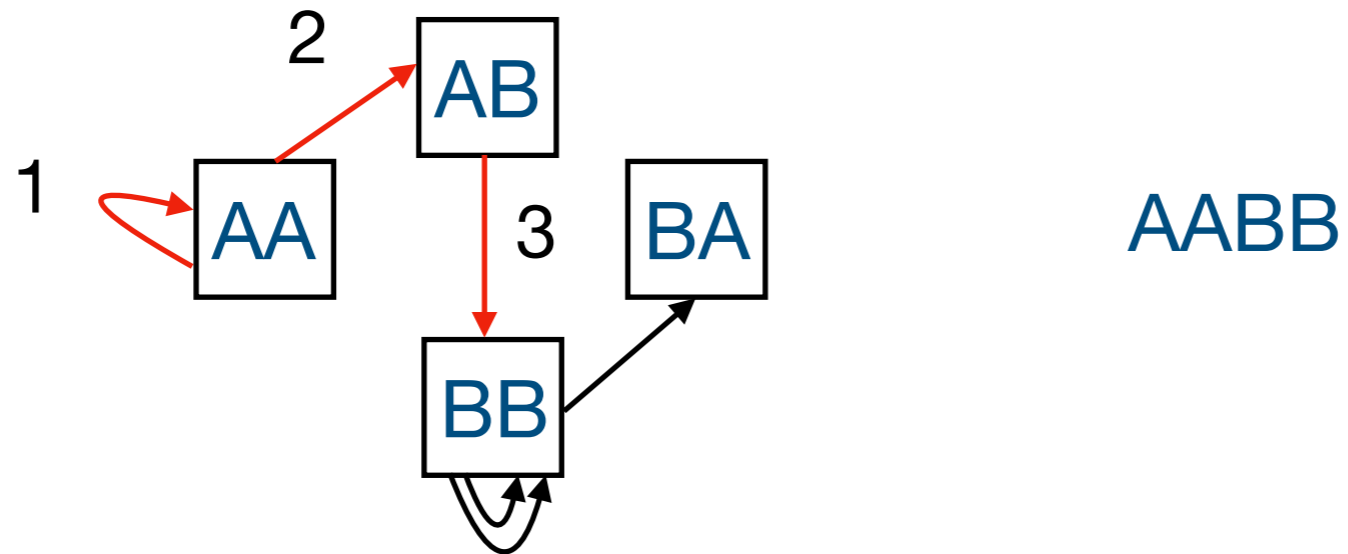
- Walk crossing each edge exactly once gives a reconstruction of the genome.

# Eulerian walk



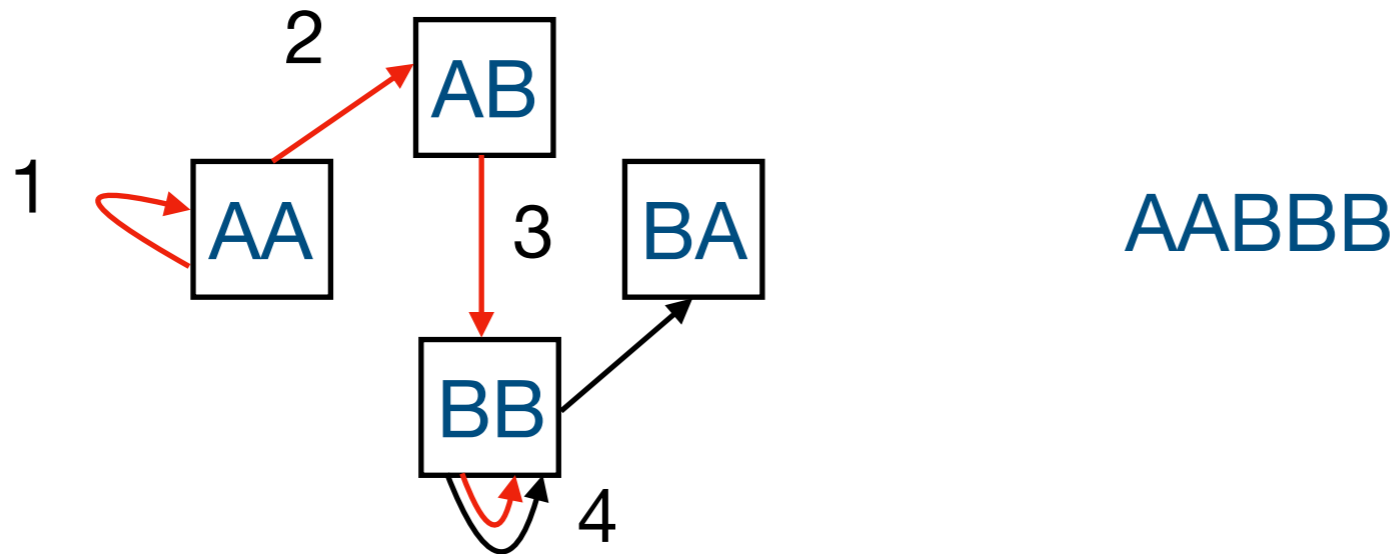
- Walk crossing each edge exactly once gives a reconstruction of the genome.

# Eulerian walk



- Walk crossing each edge exactly once gives a reconstruction of the genome.

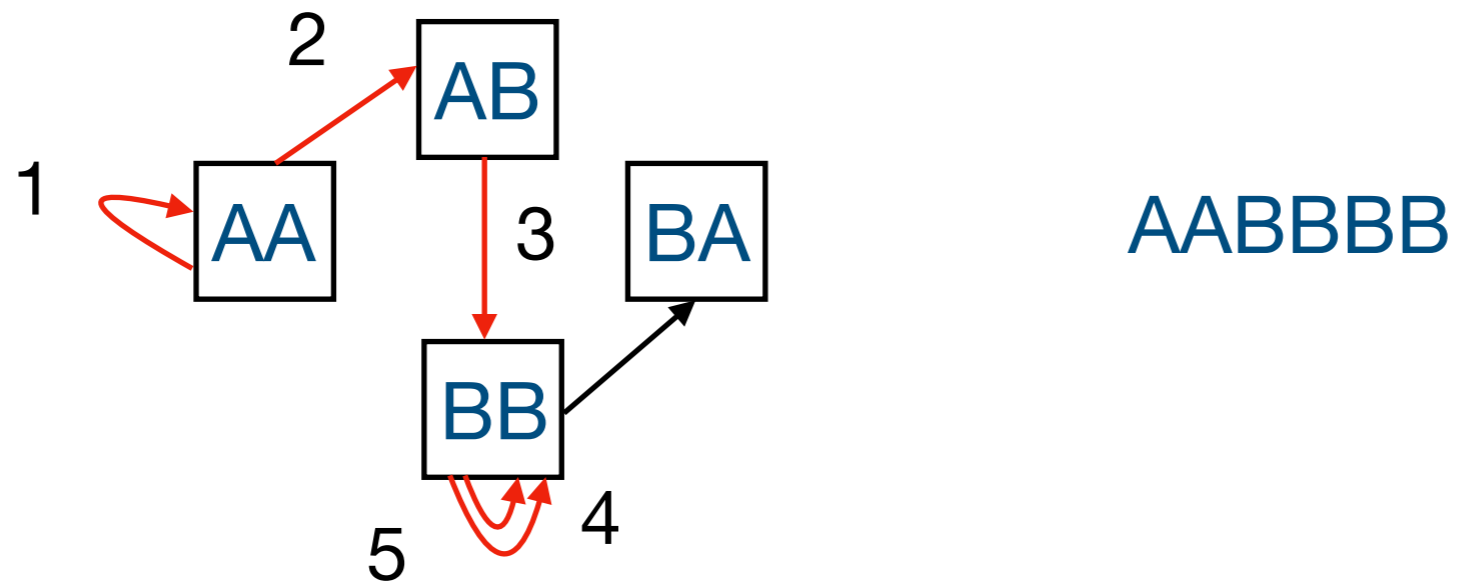
# Eulerian walk



AABB

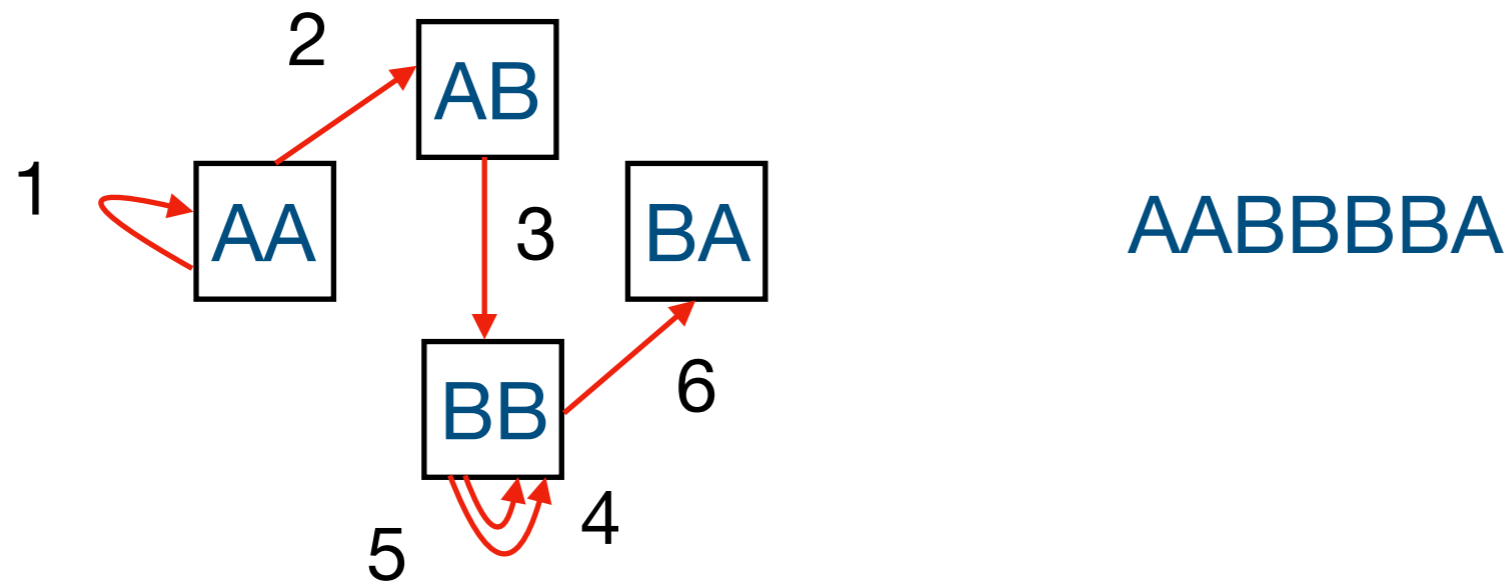
- Walk crossing each edge exactly once gives a reconstruction of the genome.

# Eulerian walk



- Walk crossing each edge exactly once gives a reconstruction of the genome.

# Eulerian walk



- Walk crossing each edge exactly once gives a reconstruction of the genome.

# SPAdes

Assembly	NG50	# contigs	Largest	Total length	MA	MM	IND	GF (%)	# genes
<b>Single-cell <i>E. coli</i></b>									
A5	14399	745	101584	4441145	3	11.92	0.19	89.867	3443
ABYSS	68534	<b>179</b>	178720	4345617	6	3.49	0.83	88.265	3704
CLC	32506	503	113285	4656964	<b>1</b>	5.54	1.00	92.286	3767
EULER-SR	26662	429	140518	4248713	12	9.98	20.17	84.846	3410
Ray	45448	361	210820	4379139	16	5.29	1.24	88.345	3634
SOAPdenovo	1540	1166	51517	2958144	<b>1</b>	<b>1.49</b>	<b>0.11</b>	57.668	1766
Velvet	22648	261	132865	3501984	2	2.19	1.17	73.761	3079
E+V-SC	32051	344	132865	4540286	2	2.26	0.70	91.727	3767
IDBA-UD contigs	98306	244	<b>284464</b>	<b>4814043</b>	3	4.37	0.23	<b>95.158</b>	<b>4041</b>
IDBA-UD scaffolds	109057	229	<b>284464</b>	<b>4813609</b>	3	4.42	0.75	<b>95.145</b>	<b>4046</b>
SPAdes 3.12 contigs	105885	231	268283	<b>4795250</b>	3	2.02	0.30	<b>94.853</b>	<b>4028</b>
SPAdes 3.12 scaffolds	<b>117600</b>	214	<b>285212</b>	<b>4800301</b>	3	2.41	0.61	<b>94.886</b>	<b>4030</b>

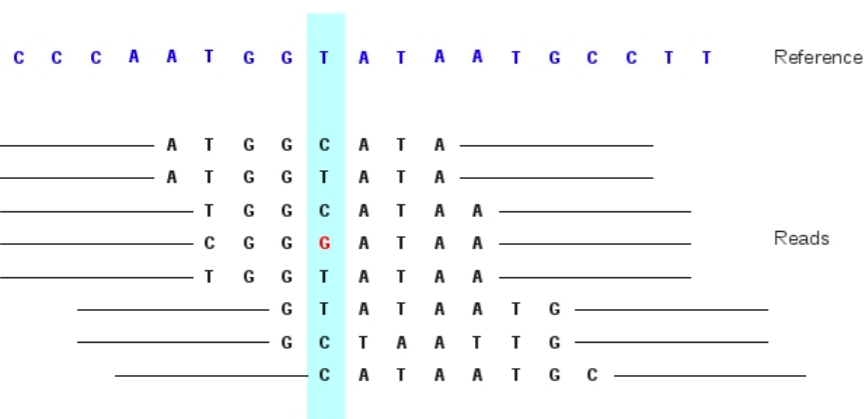
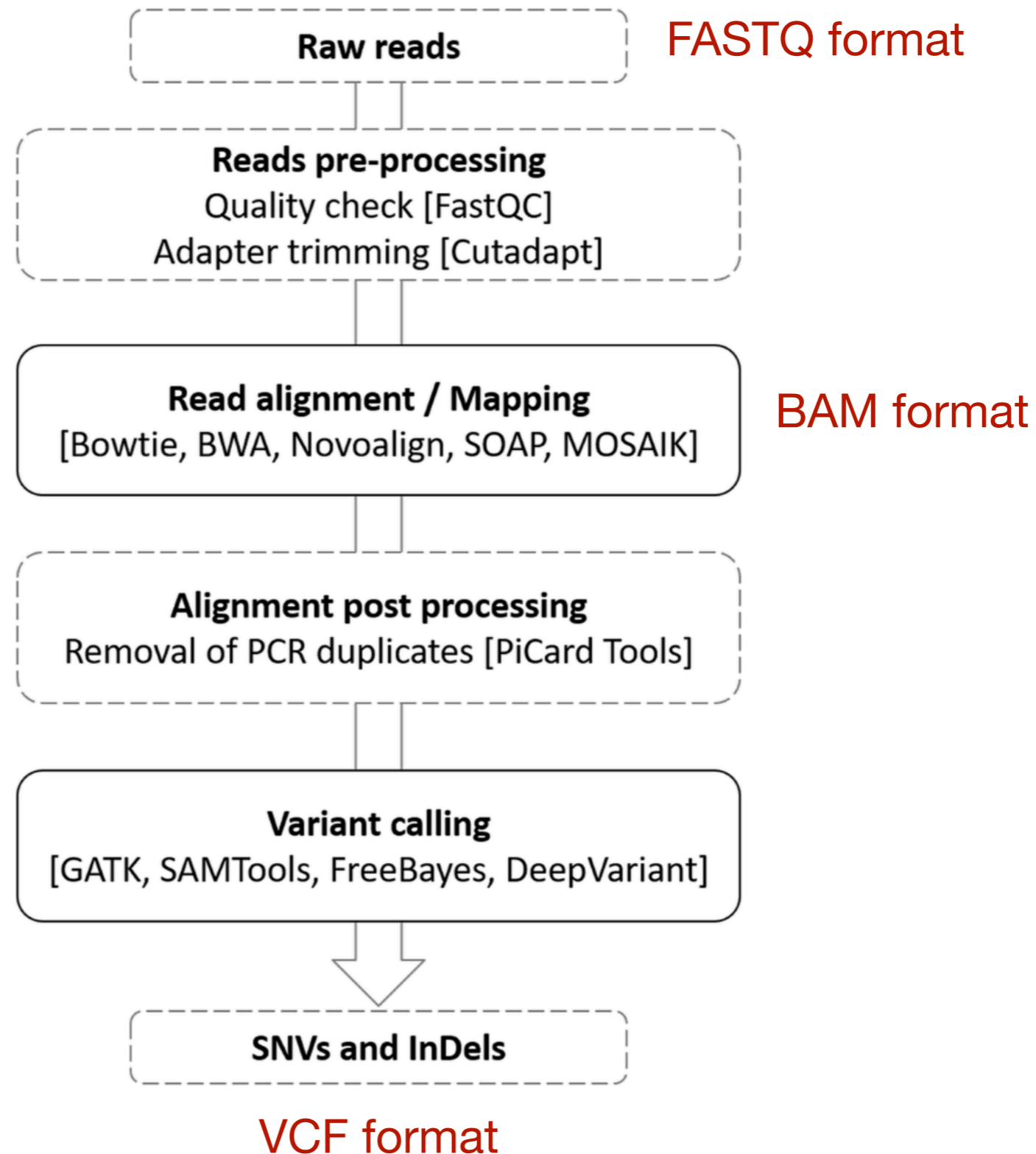
<https://cab.spbu.ru/software/spades/#examples>

- SPAdes is a de-brujin graph based genome assembler.
- By default SPAdes assembles using kmers of lengths 21, 33, and 55 and chooses the assembly with the best N50 score.
- N50 can be understood as the median contig length in the assembly.



# **Detection of genetic variants**

# Variant calling pipeline



Kumaran, M., Subramanian, U., & Devarajan, B. (2019). Performance assessment of variant calling pipelines using human whole exome sequencing and simulated data. *BMC bioinformatics*, 20(1), 1-11.

# VCF format

The diagram illustrates the VCF format structure, divided into a **VCF header** and a **Body**.

**VCF header:** Contains mandatory and optional header lines. Mandatory lines include file format, date, source, reference, and format definitions. Optional lines provide meta-data about annotations.

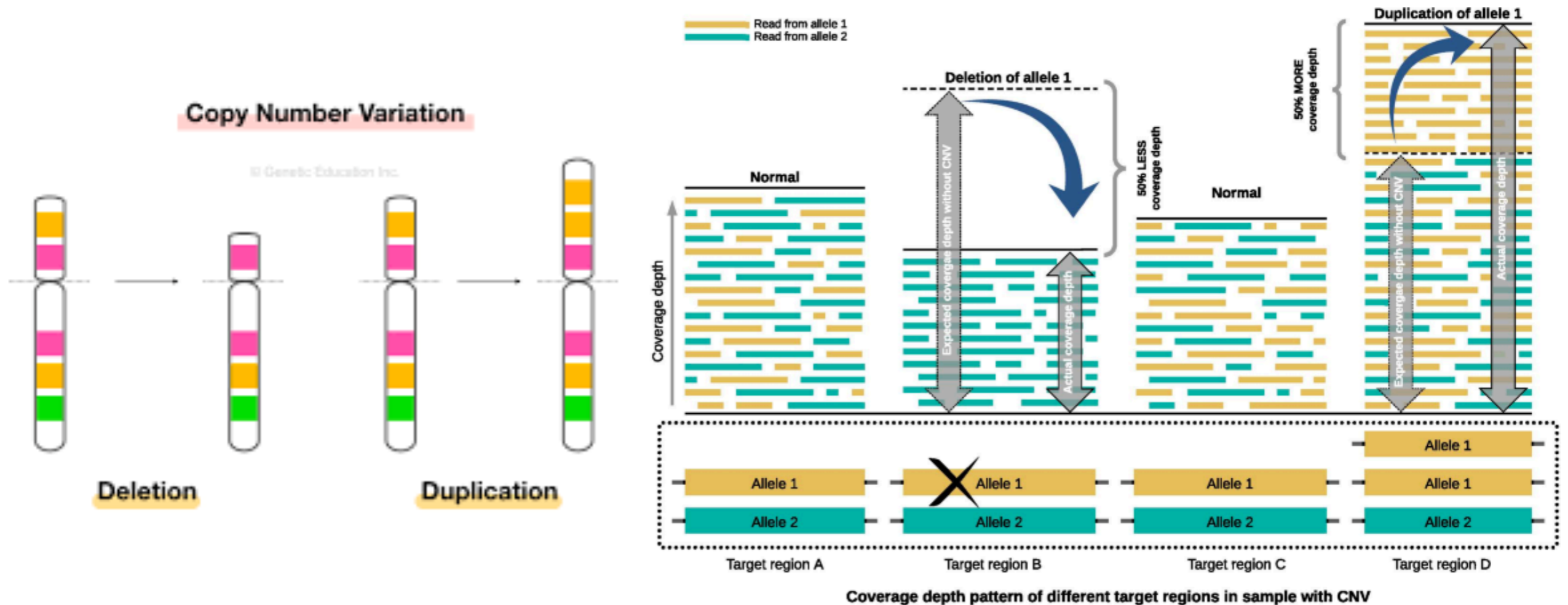
**Body:** A table of variant calls with columns: CHROM, POS, ID, REF, ALT, QUAL, FILTER, INFO, FORMAT, SAMPLE1, and SAMPLE2. Annotations identify variant types like Deletion, SNP, Large SV, and Insertion, as well as phased data and alternate alleles.

```

##fileformat=VCFv4.0
##fileDate=20100707
##source=VCFtools
##reference=NCBI36
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality (phred score)">
##FORMAT=<ID=GL,Number=3,Type=Float,Description="Likelihoods for RR,RA,AA genotypes (R=ref,A=alt)">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##ALT=<ID=DEL,Description="Deletion">
##INFO=<ID=SVTYPE,Number=1,Type=String,Description="Type of structural variant">
##INFO=<ID=END,Number=1,Type=Integer,Description="End position of the variant">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT SAMPLE1 SAMPLE2
1 1 . ACG A,AT . PASS . GT:DP 1/2:13 0/0:29
1 2 rs1 C T,CT . PASS H2;AA=T GT:GQ 0|1:100 2/2:70
1 5 . A G . PASS . GT:GQ 1|0:77 1/1:95
1 100 T <DEL> . PASS SVTYPE=DEL;END=300 GT:GQ:DP 1/1:12:3 0/0:20
  
```

- The **Variant Call Format (VCF)** is the standard file format for storing genetic variation and was developed as part of the 1000 Genomes Project.

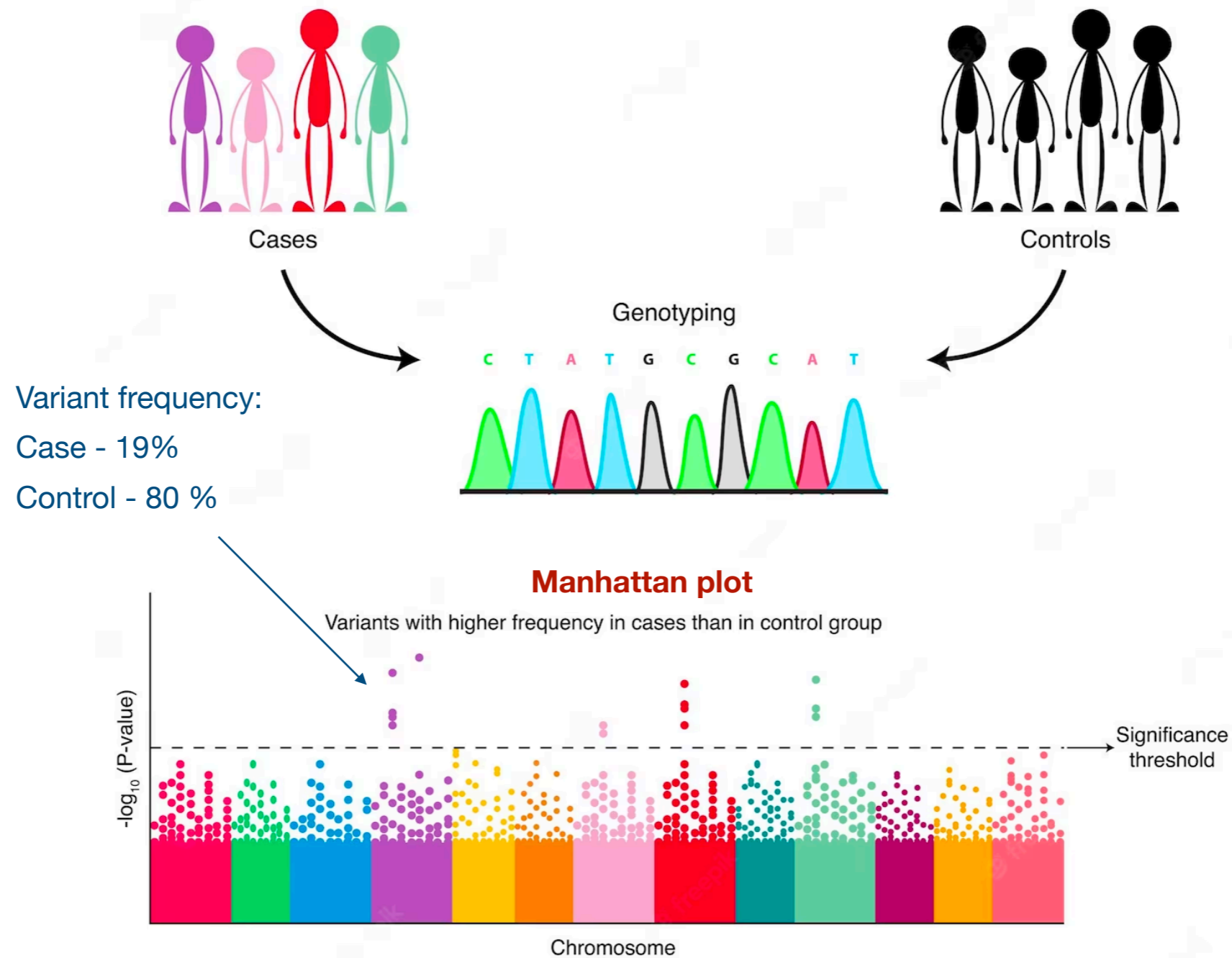
# Copy number variation detection



- CNVs are regions of the genome with variable number of copies.
- DNA-Seq can detect CNVs by analyzing the number of sequencing reads that map to a genomic region.
- Higher reads suggest a duplication, while lower reads suggest a deletion.
- CNV detection requires careful normalization and calibration, as read depth can be affected by factors such as GC content, mappability, and sequencing bias.

# **Genome-wide association analysis**

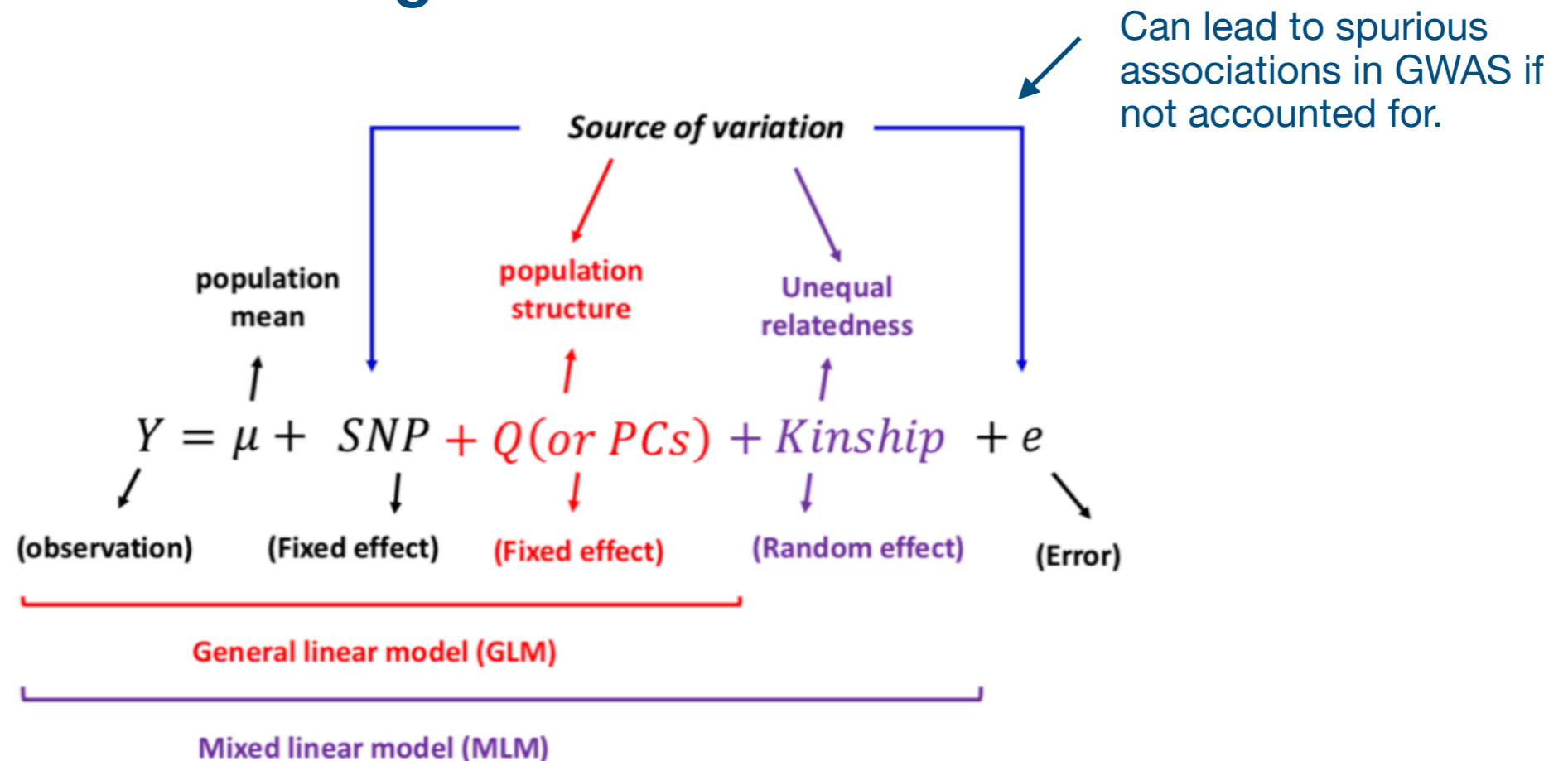
# GWAS at a glance



**GWAS** aim to identify if any of the millions of SNPs are associated with a specific disease (e.g. Cancer) or trait (e.g. Height / intelligence).

# GWAS test using linear regression

## Account for confounding variables



### Population Structure:

- Systematic differences in allele frequencies between subgroups in a population due to non-random mating between individuals.
- Can be estimated from data using statistical methods such as PCA.

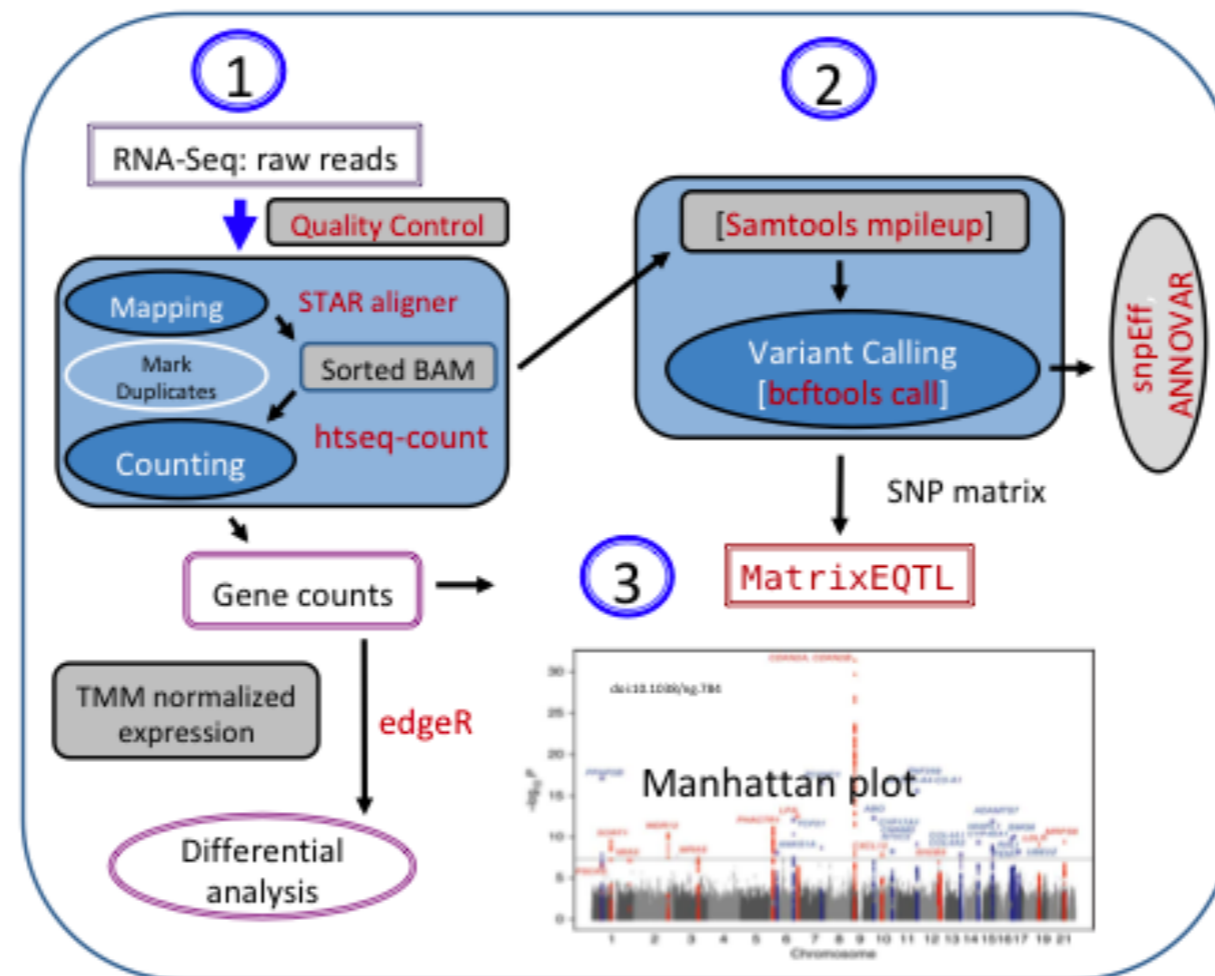
### Kinship:

- Describes the genetic relatedness between individuals in a population.
- Kinship matrix is often modeled as the covariance of the random effect term.

# **Expression quantitative trait loci**



# General workflow



<https://adinasarapu.github.io/posts/2017/12/blog-post-eqtl/>

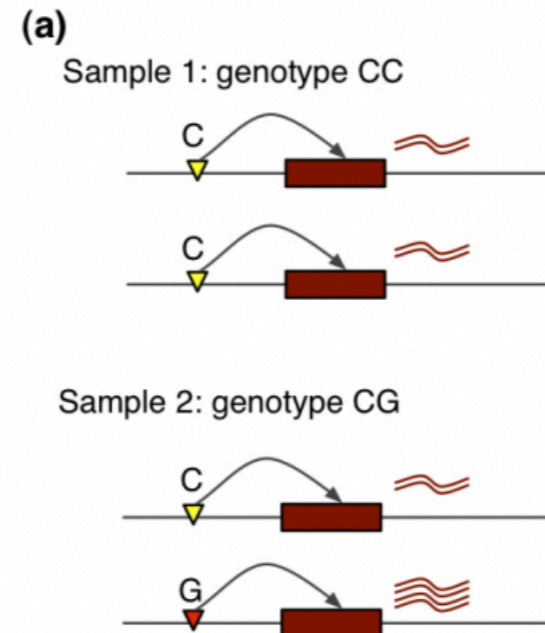
- An expression quantitative trait locus (**eQTL**) is a genetic locus that affects gene expression.
- eQTL mapping studies use RNA-Seq data to identify eQTLs.
- Variants are called either from DNA-Seq / RNA-Seq; expression levels are quantified via the regular pipeline, and differential analysis is performed between genotypes.

# Cis & trans QTL

## cis-eQTL:

- Variants affecting expression of **local genes**.
- Found in promoter and gene body of the effected genes.

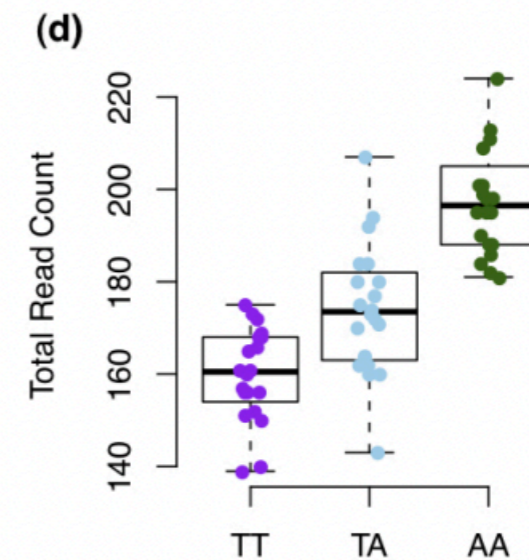
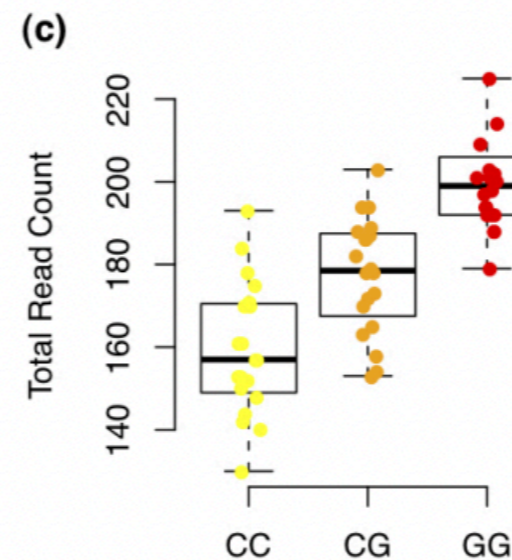
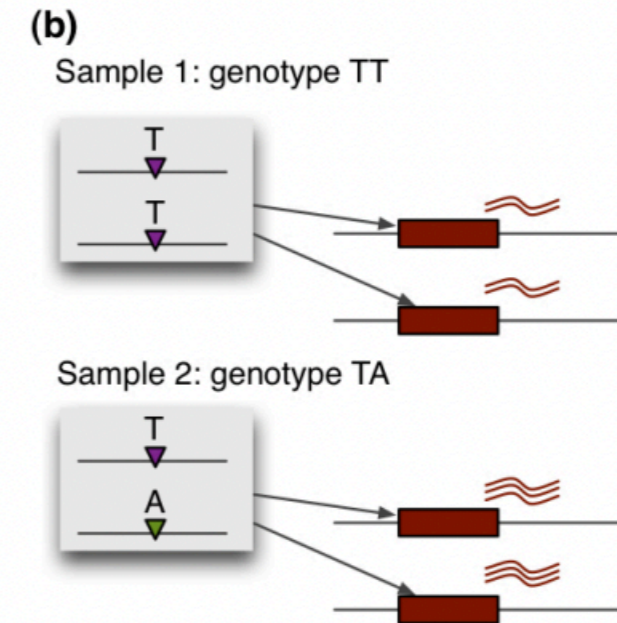
## cis-eQTL



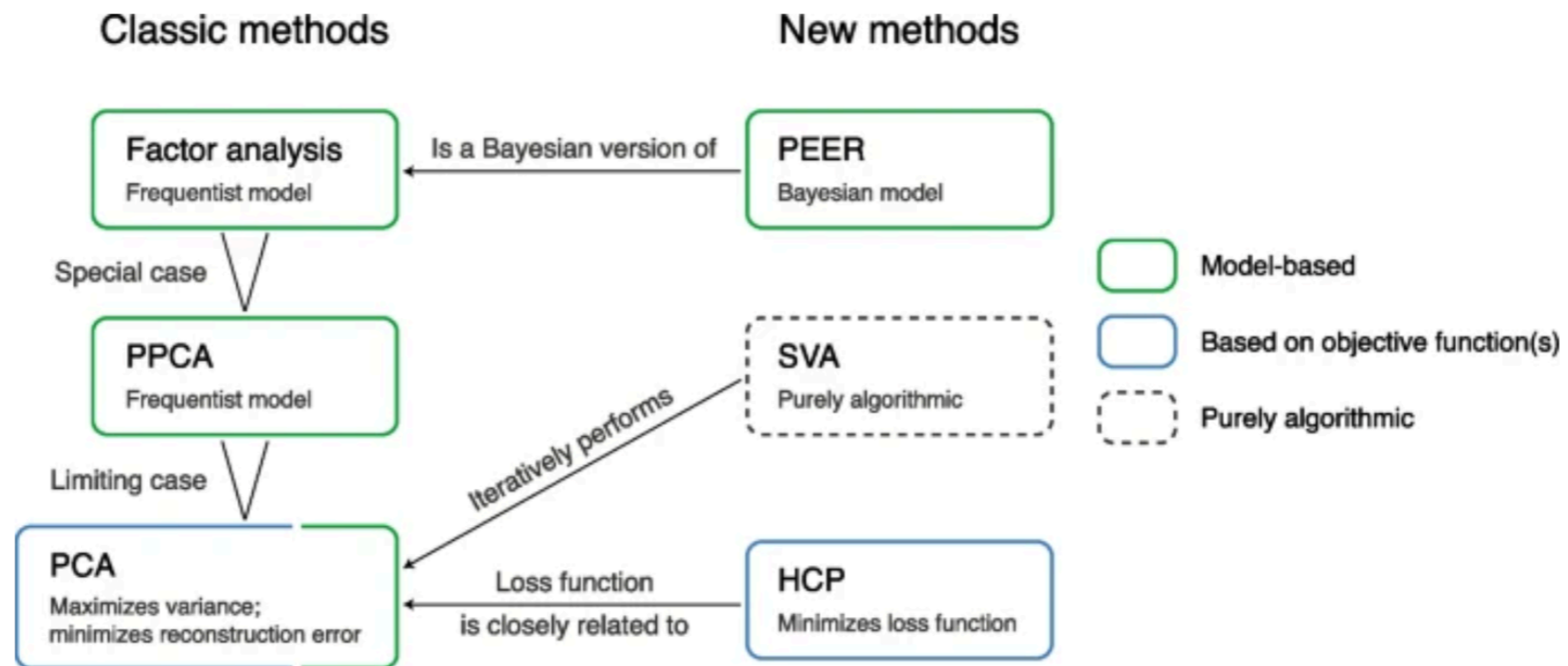
## trans-eQTL:

- Variants affecting expression of **distal genes**.
- Found in other regions.

## trans-eQTL



# Accounting for hidden batches



- P values of eQTL association are calculated from the linear regression tests with the following regression equation.
- $$\text{expression} = \alpha + \sum_k \beta_k \cdot \text{covariate}_k + \gamma \cdot \text{genotype}$$
- The  $K$  covariates are estimated by hidden variable inference methods such as PCA, SVA, PEER or HCP.

# Accounting for hidden batches

PCA  
outperform  
the rests ...



- P values of eQTL association are calculated from the linear regression tests with the following regression equation.
- $\text{expression} = \alpha + \sum_k \beta_k \cdot \text{covariate}_k + \gamma \cdot \text{genotype}$
- The K covariates are estimated by hidden variable inference methods such as PCA, SVA, PEER or HCP.

# **In silico mutation analysis**

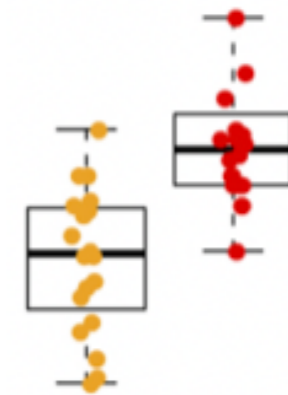
# Infer the consequence of mutation by predictive modeling



**GWAS:** Variants  $\langle - - \rangle$  Phenotype

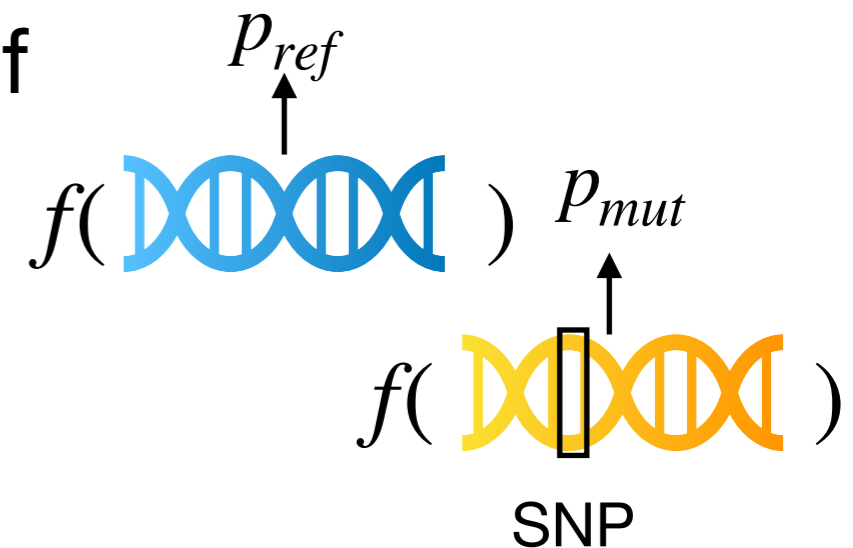
Only correlation !

**eQTL:** Variants  $\langle - - \rangle$  Gene expression



**In-silico mutation:** Variants  $\langle - - \rangle$  Change of Prob(epigenetic marker)

Can reveal causality !



# In-silico mutation

- $f()$  is a sequence based predictive model, it accepts an input of a DNA string and output a probability of the string being a functional epigenetic modification or protein.
- Calculate the probabilities of WT sequence and mutated sequence (e.g. caused by a SNP)

$$f(\text{WT sequence}) \rightarrow \text{prob1}$$

$$f(\text{mutated sequence}) \rightarrow \text{prob2}$$

- Inference of SNP function:
  - $\text{prob1} \gg \text{prob2}$ : **loss of function** mutation
  - $\text{prob1} \ll \text{prob2}$ : **gain of function** mutation