

review

Alignment algorithms

Genome
sequences



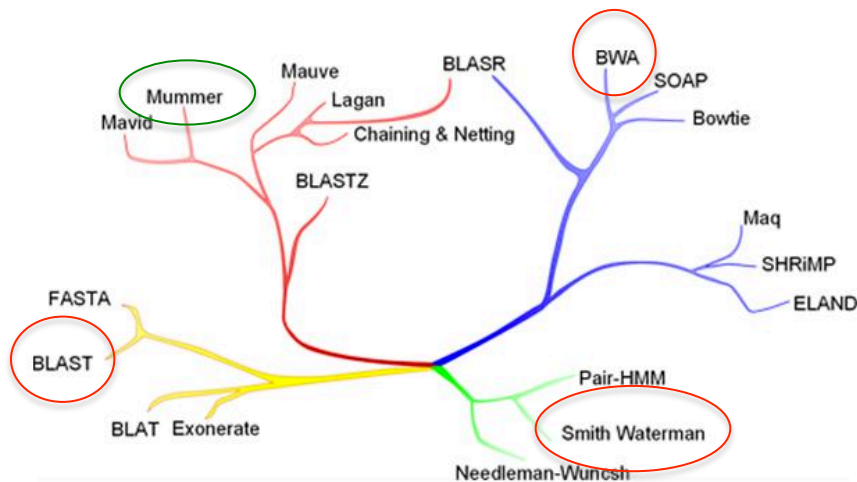
“sorted or
indexed”
genome
memory

one by one

reads

alignments

Aligner phylogeny



Whole genome
Pairwise heuristic

Short read
Sensitive global aligners

SAM and BAM alignment format

- **SAM** stands for Sequence Alignment/Map format that is a generic alignment format for storing read alignments against reference sequence.
- The **BAM** format, the binary representation of SAM, contains exactly the same information as SAM,
- The SAM/BAM, together with **SAMtools**, separates the alignment step from downstream analyses, enabling a generic and modular approach to the analysis of genomic sequencing data.

Table. file size of two example files

File type	Storage usage
SAM	313 M
BAM	97 M

Outline

- Overview of genomic variants
- Data for variant discovery
- Bioinformatics of variant discovery
- the methods for variant (SNP) validation

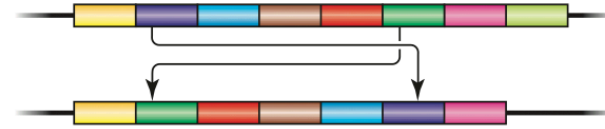
Genomic variants (ploymorphism)

1. SNP

Point mutation

TGCATT **G** CGTAGGC
 ↓
TGCATT **C** CGTAGGC

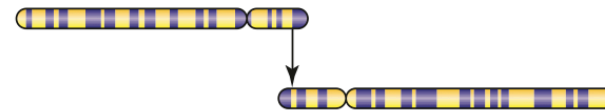
Inversion



Insertion

TGCATTTAGGC
TGCATT **CCG** TAGGC
 ↑

Chromosome fusion

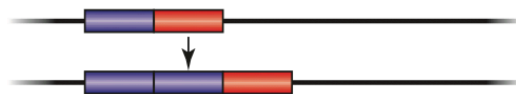


2. INDEL

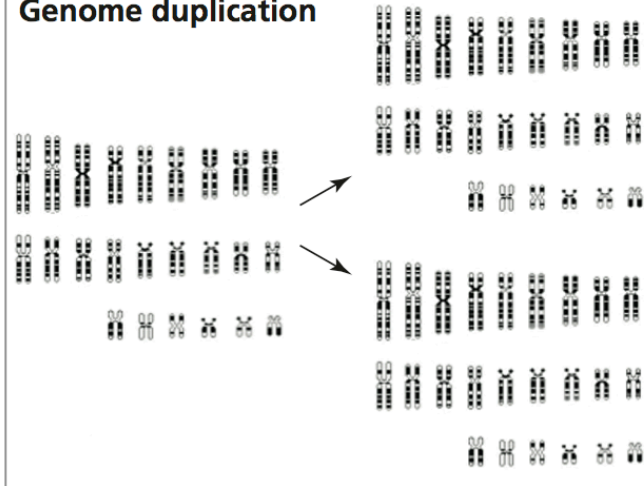
Deletion

TGCATT ~~**CCG**~~ TAGGC
 ↓
TGCATT TAGGC

Gene duplication



Genome duplication



3. genomic structural variation

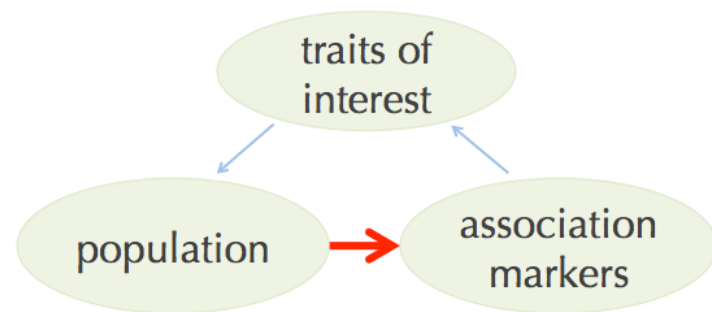
- copy number variation (presence-absence variation)
- other re-arrangements

Genomic variants - SNPs

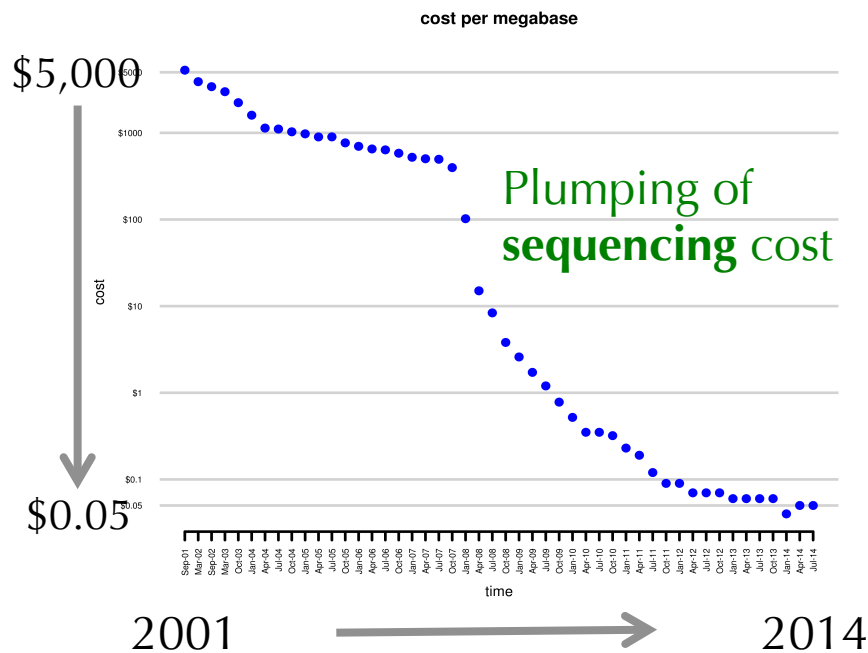
- SNP is a single nucleotide polymorphism.
- Frequencies of SNPs are depended on species. For example, millions of SNPs have been discovered in human.
- Most SNPs are bi-allelic. (mutation rate per site is about 10^{-8})
- Most have no effects on cell function but some could have important phenotypic consequence.

Applications of SNPs

1. Genetic markers to map the genetic controlling of traits (quality traits, quantitative traits, gene expression, etc)
2. Genetic markers to construct genetic maps
3. Markers to construct phylogenetic trees



Next-Generation Sequencing to generate data for variant discovery



GATCTGCGT**C**ATACGGAAT
GATCTGCGT**G**ATACGGAAT
GATCTGCGT**C**ATACGGAAT
GATCTGCGT**G**ATACGGAAT
GATCTGCGT**G**ATACGGAAT
GATCTGCGT**G**ATACGGAAT
GATCTGCGT**C**ATACGGAAT
GATCTGCGT**C**ATACGGAAT
GATCTGCGT**G**ATACGGAAT
GATCTGCGT**C**ATACGGAAT

-----**C/G**-----

Approaches for data generation

- **Whole genome sequencing (WGS):** high genome coverage but costly for large genomes
- **Exome-capture sequencing:** target on genic regions but still expensive to perform large number of samples
- **RNA sequencing (RNA-Seq):** obtain data on genic regions and provide expression information
- **Genotyping-By-Sequencing (GBS):** cost-efficient and high-throughput approach

Data to variant (SNP) discovery

- Data: sequencing reads
- Reference-based approach
 1. Alignment-based SNP discovery (standard)
 2. Assembly-based SNP discovery
- Reference-free approach

Usually used in the sequencing projects with sequencing data from multiple individuals

Alignment-based SNP discovery

...GATCTGCGTCATACGGAAT... (reference)

GATCTGCGT**G**ATACGGAAT
GATCTGCGT**C**ATACGGAAT
GATCTGCGT**G**ATACGGAAT
GATCTGCGT**G**ATACGGAAT
GATCTGCGT**G**ATACGGAAT
GATCTGCGT**C**ATACGGAAT
GATCTGCGT**C**ATACGGAAT
GATCTGCGT**G**ATACGGAAT
GATCTGCGT**C**ATACGGAAT

} reads

-----C/G-----

Alignment-based SNP discovery, cont.

General procedure

- Reads cleanup (adaptor, quality trimming, e.g., trimmomatic)
- Reads aligned to the reference genome with aligners
 1. BWA, Bowtie (DNA-Seq reads)
 2. GSNAP, Tophat (RNA-Seq reads)
- Post-alignment filtering and convert SAM (alignment file) to BAM
- SNP calling with software packages: Samtools, GATK, VarScan2
- Use population information or some criteria to filter SNP sets

an example of the alignment of a RNA-Seq read using GSNAP

[illegible]

Interpretation of the BWA alignment

SAM output:

```
HWI-ST897:104:C015GACXX:6:1101:12678:20443 163 U00096
1888286 60 64M1D20M = 1888358 170
GCCAACAGCCGCGACTTCCTGTACGCCAGGATGCTGCATGACGACATCTTCAATCTCGTTGGGAAGACGTTAAAAACGGAAACC CCCFFFFFHHFHHJJJJJJJ
JHIJHIJIIJIIJJJJJJJJJJJJJJHHFFFFFFFEEEEEEEDDDDDDA5,53,8<?CC(50?8BD3? NM:i:2 AS:i:72 XS:i:
0 RG:Z:S1
```

CIGAR: 64M1D20M

NM: edit distance

edit distance is a way to quantify the dissimilarity of two strings (e.g., words) by counting the minimum number of edits (substitution, insertion, and deletion) required to transform one string into the other.

fact -> fit (2)

AACCT -> AAAC (1)

AACCT -> ACCTA (?)

Polymorphism based on Alignment + reference genome

SAM output (BWA):

```
HWI-ST897:104:C015GACXX:6:1101:12678:20443 163 U00096
1888286 60 64M1D20M = 1888358 170
GCCAACAGCCGCGACTTCCTGTACGCCAGGATGCTGCATGACGACATCTTCAATCTCGTTGGGAAGACGTTAAAAACGGAAACC CCCFFFFFFHHFHHJJJJJJJ
JHIJHIJIIJIIJJJJIIJJJJIIJHHFFFFFFFEEEEEDDDDDDA5,53,8<?CC(50?8BD3? NM:i:2 AS:i:72 XS:i:
0 RG:Z:S1
```

mapping position and CIGAR determine the alignment

```
Query 1 GCCAACAGCCGCGACTTCCTGTACGCCAGGATGCTGCATGACGACATCTTCAATCTCGTT 60
|||||
Sbjct 1888286 GCCAACAGCCGCGACTTCCTGTACGCCAGGATGCTGCATGACGACATCTTCAATCTCGTT 1888345

Query 61 GGGA-AGACGTTAAAAACGGAAACC 84
||| |||||
Sbjct 1888346 GGGATAGACGTTAAAACCGAAACC 1888370
```

Alignment-based SNP discovery: GATK (1)

- The Genome Analysis Toolkit (GATK) is a software package developed at the Broad Institute to primarily focus on variant discovery and genotyping.
- Input data: BAM files and reference genome
- Required tools: Picard and Samtools
- Code example:

```
java -jar GenomeAnalysisTK.jar \  
    -T UnifiedGenotyper \  
    -R your_reference \  
    -I your_bam \  
    -glm BOTH
```

```
### BOTH = SNP + INDEL
```

GATK (2)

VCF (Variant Call Format) output

<https://samtools.github.io/hts-specs/VCFv4.2.pdf>

isolate 1



isolate 2



#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	DH10B	MG1655
ref1	89089	.	C	A	782.76	GT:AD:DP:GQ:MLPSAC:MLPSAF:PL	1:0,18:18:99:1:1.00:781,0	0:27,0:27:99:0:0.00:0,1149
ref1	89103	.	G	C	690.76	GT:AD:DP:GQ:MLPSAC:MLPSAF:PL	1:0,16:16:99:1:1.00:689,0	0:29,0:29:99:0:0.00:0,1253
ref1	89143	.	A	G	448.76	GT:AD:DP:GQ:MLPSAC:MLPSAF:PL	1:0,11:11:99:1:1.00:447,0	0:27,0:27:99:0:0.00:0,1165
ref1	89145	.	G	T	405.76	GT:AD:DP:GQ:MLPSAC:MLPSAF:PL	1:0,10:10:99:1:1.00:404,0	0:28,0:28:99:0:0.00:0,1215

GT: AD : DP: GQ: MLPSAC: MLPSAF: PL

1 : 0,18: 18: 99: 1 : 1.00 : 781,0

GT=Genotype (0 or 1)

AD=Allelic depths for the ref and alt alleles

DP=Approximate read depth

GQ=Genotype Quality

MLPSAC=Maximum likelihood expectation (MLE) for the alternate allele count

MLPSAF=Maximum likelihood expectation (MLE) for the alternate allele fraction

PL=Normalized, Phred-scaled likelihoods for genotypes

$$\text{Prob}(0) = 10^{(-781/10)} = 7.9\text{e-}79 \quad \text{Prob}(1) = 10^{(-0/10)} = 1$$

GATK (3)

- GATK can be used to filter SNPs.

```
java GenomeAnalysisTK.jar \  
  -T SelectVariants \  
  -R your_reference \  
  --variant your_vcf \  
  -select 'DP >= 3.0' \  
  --restrictAllelesTo BIALLELIC \  
  --selectTypeToInclude SNP
```

- Filter variants based on the experimental purpose and genetic features

Falsely discovered SNPs

Can you think about what could result in falsely discovered SNPs using alignment-based SNP methods?

Alignment-based SNP discovery: alignment issues

- Misalignments
- Genome duplications
- Highly divergent regions

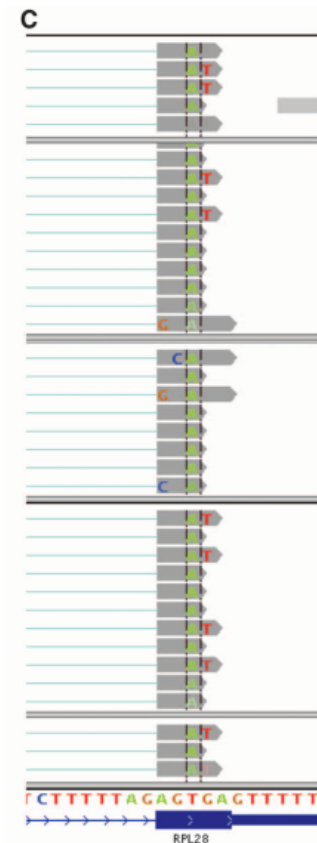
Examples:



Widespread RNA and DNA Sequence Differences in the Human Transcriptome
Mingyao Li *et al.*
Science **333**, 53 (2011);
DOI: 10.1126/science.1207018

The misalignments of RNA-Seq data or DNA-Seq data led to this discovery

Comment on "Widespread RNA and DNA Sequence Differences in the Human Transcriptome"
Claudia L. Kleinman and Jacek Majewski
Science **335**, 1302 (2012);
DOI: 10.1126/science.1209658



Assembly-based SNP discovery

- Cortex (Iqbal *et al.*, 2012 Nature Genetics)

de novo assembly and graphic comparison for variant discovery

- Fermi (Li H, 2012 Bioinformatics)

de novo assembly to unitigs* and then alignment to the reference genome for variant discovery

(Conceptually, unitigs are confident contigs)

- Discover (Neil *et al.*, 2014 Nature Genetics)

Region *de novo* assembly to contigs and then alignment to the reference genome for variant discovery

Table 2 Estimated sensitivity and specificity of variant call sets

Call set	Read length (bp)	Percent false negatives	Number of heterozygous/homozygous variants	Percent false positives		
				Heterozygous variants	Homozygous variants	All variants
GATK-250	250	12.3 ± 1.8	1.54	1.82 ± 0.45	0.74 ± 0.72	1.39 ± 0.39
Cortex-250	250	39.3 ± 2.6	1.39	0.33 ± 0.18	3.46 ± 0.61	1.64 ± 0.28
DISCOVAR-250	250	06.0 ± 1.2	1.57	1.44 ± 0.23	1.94 ± 0.40	1.63 ± 0.21

Variant annotation

Gene coding regions

- *Synonymous*: changes that do not alter the encoded amino acid
- *Non-synonymous*
 1. *Missense*: changes that alter encoded amino acid
 2. *Nonsense*: changes that produce a stop codon from an amino acid codon, resulting in a shortened protein
- *Frameshift* (caused by insertion/deletion)

Splicing sites

Of an intron, a donor site (5' end of the intron) and an acceptor site (3' end of the intron) are required for splicing.

Variant annotation - SnpEff

SnpEff is a variant annotation and effect prediction tool. It annotates and predicts the effects of variants on genes.

Input data:

- Genome annotation database
- Variant data: VCF file

Running:

```
java -jar snpEff.jar GRCh37.75 my.vcf
```

Effect
INTERGENIC
UPSTREAM
UTR_5_PRIME
UTR_5_DELETED
START_GAINED
SPLICE_SITE_ACCEPTOR
SPLICE_SITE_DONOR
START_LOST
SYNONYMOUS_START
CDS
GENE
TRANSCRIPT
EXON
EXON_DELETED
NON_SYNONYMOUS_CODING
SYNONYMOUS_CODING

Detailed effect list from SnpEff

FRAME_SHIFT
CODON_CHANGE
CODON_INSERTION
CODON_CHANGE_PLUS_CODON_INSERTION
CODON_DELETION
CODON_CHANGE_PLUS_CODON_DELETION
STOP_GAINED
SYNONYMOUS_STOP
STOP_LOST
INTRON
UTR_3_PRIME
UTR_3_DELETED
DOWNSTREAM
INTRON_CONSERVED
INTERGENIC_CONSERVED

SNP genotyping

Large-scale (thousands to approximate 1 million)

- Illumina Beadchip (beads hybridization based)
- Affymetrix SNP array (microarray-hybridization-based)

Medium-scale (hundreds of markers)

- Fluidigm (Microfluidic-based)
- Sequenom iPLEX (mass spectrometry method)

Small-scale

- High Resolution Melt (HRM (melting))
- KASP
- Taqman

SNP validation

- Cross-checking using different datasets or platforms
- Genetic mapping of SNPs
- Expectation from certain genome materials (e.g., bacterial genome and inbred lines)

Summary

- The strategy to generate data for SNP discovery is depended on experimental purpose, genetic features of the population, timetable, and budget.
- A standard approach for SNP discovery is through mapping reads to the reference sequences, thereby identifying variants between reads and reference. The most popular method is GATK.
- More flexible and cost-efficient SNP validation approaches need to be developed to leverage variant discovery.