# Whole Genome Resequencing
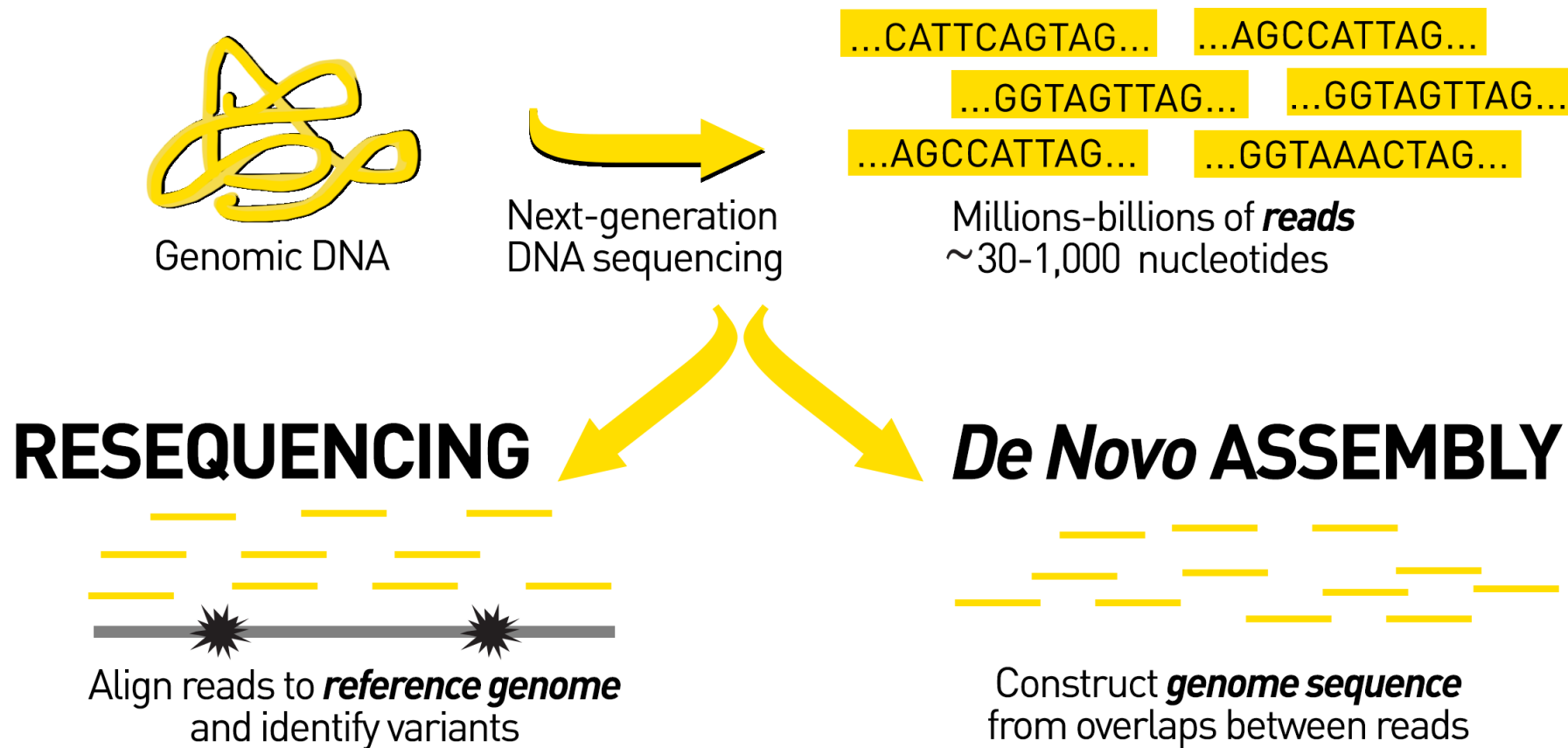
李 余 动
lyd@zjsu.edu.cn

# Topics

- **Genetic Variations**

- **Read Mapping**
  - ✦ BWA, Bowtie2…
  - ✦ SAM/Samtools

- **Variant calling**
  - ✦ Bcftools, GATK…
  - ✦ VCF format

- **SARS-CoV-2 resequencing**
  - ✦ Data analysis pipeline

# 全基因组测序(Whole Genome Sequencing)

- 当前基因组测序主要是针对已有基因组物种的重测序。
- 重测序主要比较个体基因组与参考基因组的差异



...CATTCAGTAG...   ...AGCCATTAG...
...GGTAGTTAG...   ...GGTAGTTAG...
...AGCCATTAG...   ...GGTAAACTAG...

Genomic DNA

Next-generation
DNA sequencing

Millions-billions of **reads**
~30-1,000 nucleotides

**RESEQUENCING**

Align reads to **reference genome**
and identify variants

**De Novo ASSEMBLY**

Construct **genome sequence**
from overlaps between reads

# Genetic Variations(遗传变异)

- **SNV**: Single Nucleotide Variant

- **Indel**: Insertion/Deletion

- **SV**: Structural Variants
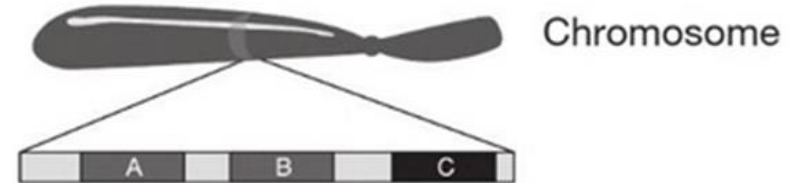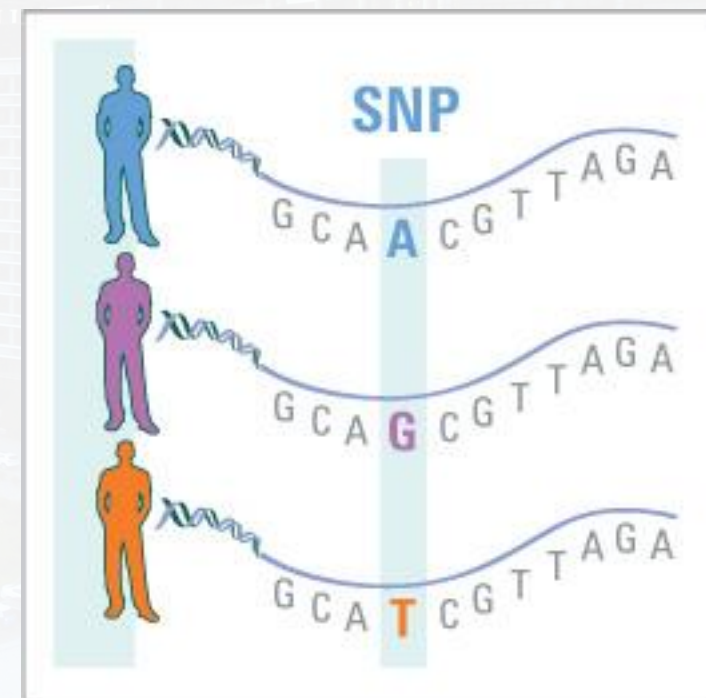
  - **CNV**: Copy Number Variation

# Single Nucleotide Polymorphisms (SNPs)

- 从种群概念上讲，一般把那些在种群中发生频率大于1%的SNVs称为SNPs。
- ~ 97 % of the genome between any two individuals is identical
  - ~ 1% of the differences are single nucleotide variations (SNPs)
  - ~2% Other changes – copy number variations, deletions
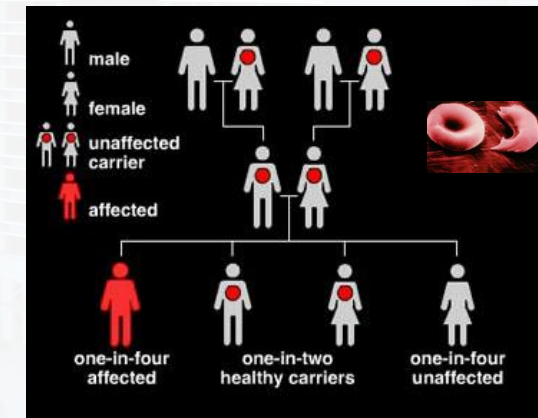- 人类基因组每隔500至1000个碱基就会存在一处SNP位点。Between 11-12 million SNPs have been identified (dbSNP)

SNP: 单核苷酸多态性

- Inherited (Germline) Mendelian genetic disorders are caused

  by variations in DNA (SNPs)

- Most of these deleterious variations affect the function of the

  encoded protein

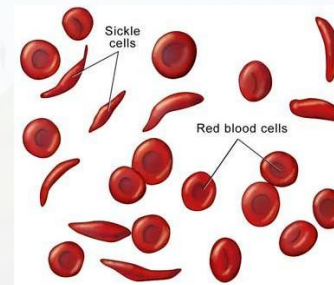  - e.g., Sickle cell anemia Val to Glu codon 6.



http://www.orgsites.com/va/pasca/

| Normal HbA | ATGGTGCACCTGACTCCTGTGGAGAAGTC |
| Disease HbS | ATGGTGCACCTGACTCCTG**A**GGAGAAGTC |

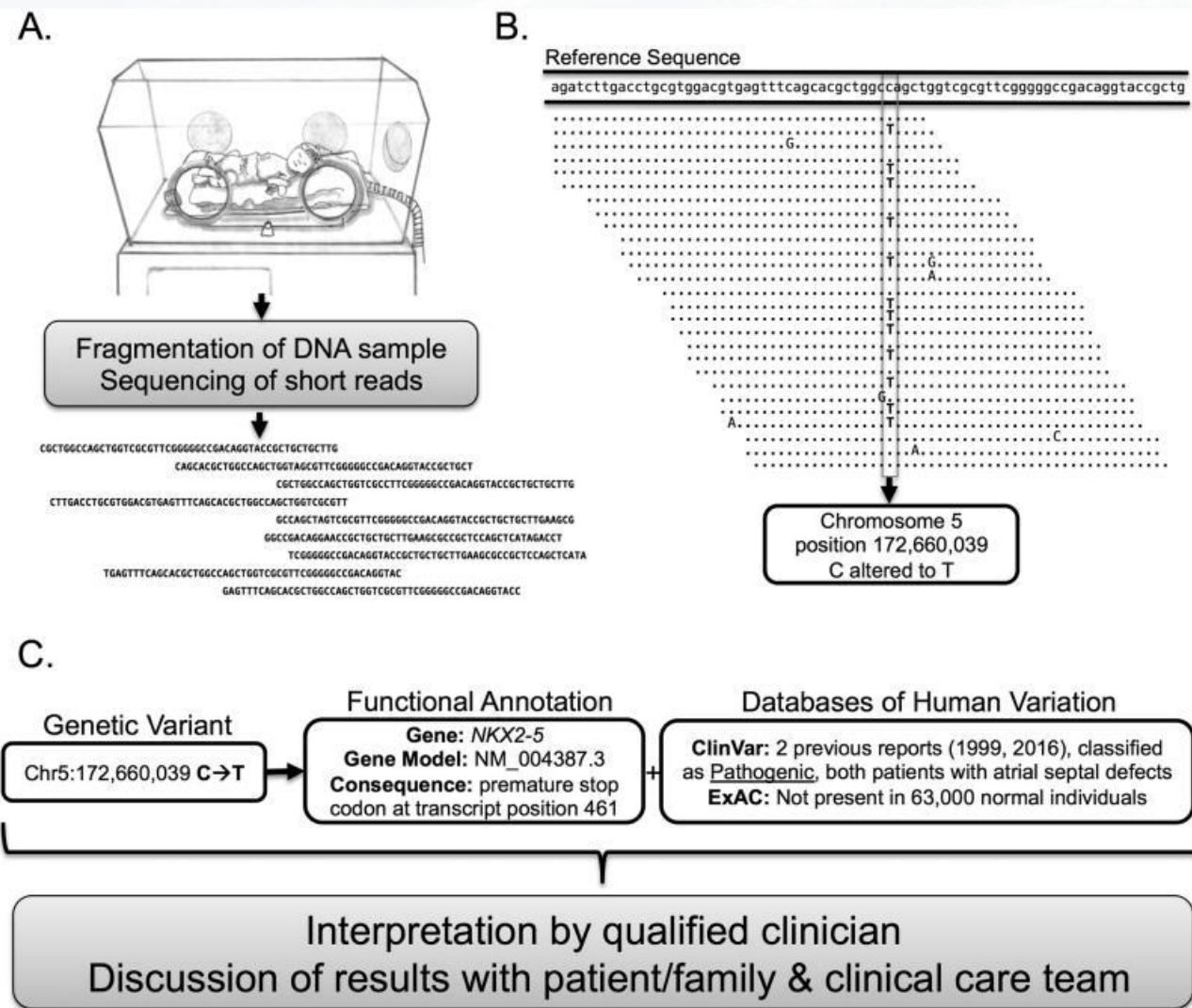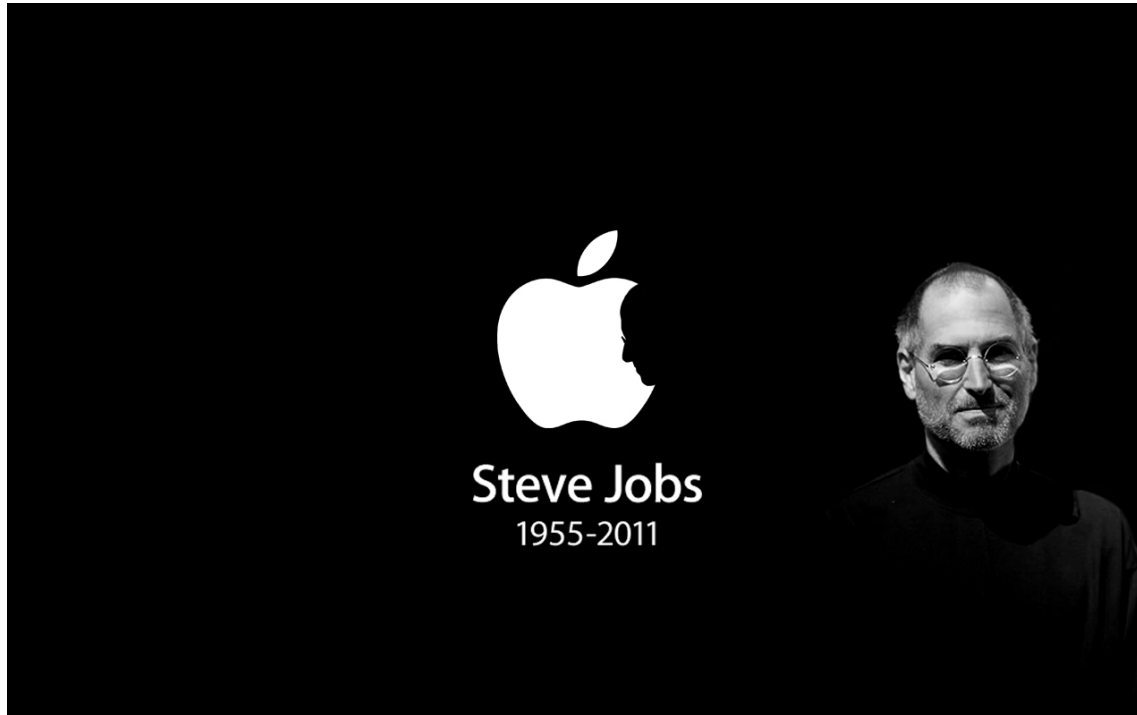| Normal HbA | MVHLTPVEKSAVTA |
| Disease HbS | MVHLTP**E**EKSAVTA |



biologycorner.com

通过不同表型差异与基因组的关联分析(Association Study)，研究表型差异的遗传学机制。

- While **single gene** and well-known Mendelian genetic disorders, such as sickle-cell anemia, Tay–Sachs disease and cystic fibrosis, can be identified with simple diagnostic techniques, **Whole genome sequencing (WGS)** can be used to identify the cause of a large range of genetic diseases.
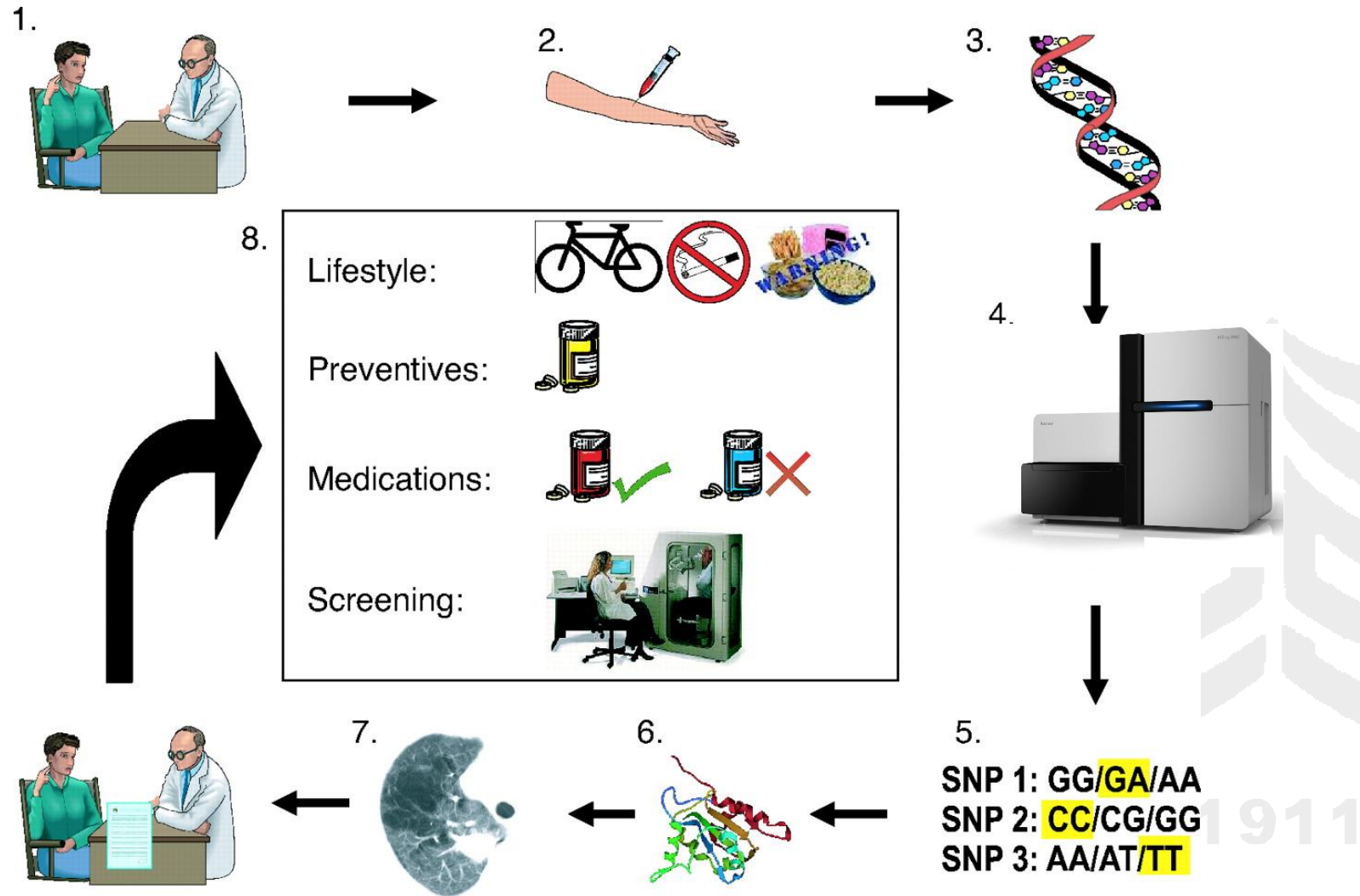
# Fighting cancer through next-generation sequencing

为了赢得与癌症的斗争，史蒂夫·乔布斯曾花费10万美元巨资为自己DNA测序，寄希望于找出治疗肿瘤的基因。

"I'm either going to be one of the first to be able to outrun a cancer like this, or I'm going to be one of the last to die from it." -- Steve Jobs

**Tebbutt S J et al. Chest 2007;131:1216-1223**
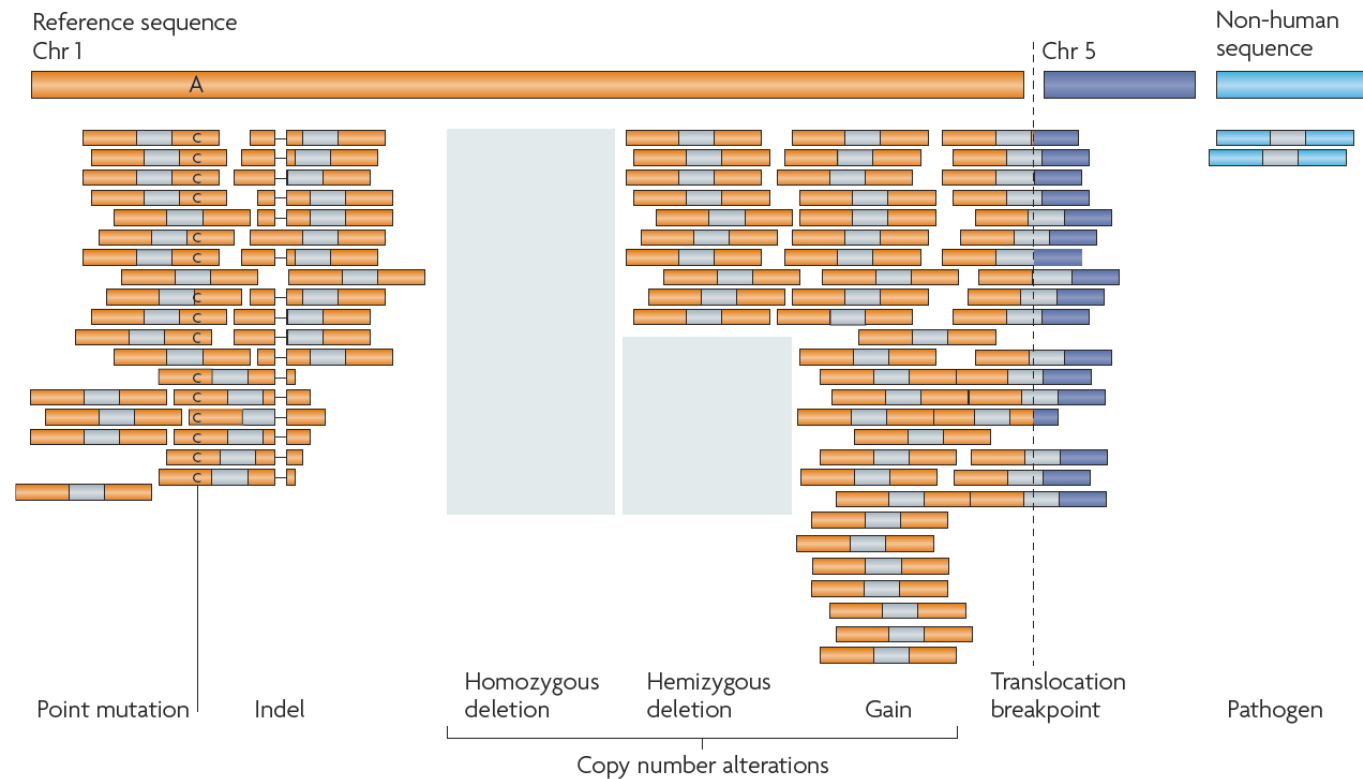
- 快速地将数百万个短读段(reads) 回帖到参考基因组上，
- 准确地鉴定SNPs和Indels等突变。



Meyerson, M., Gabriel, S. & Getz, G. Advances in understanding cancer genomes through second-generation sequencing. Nat Rev Genet 11, 685–696 (2010).

# Read Mapping : 序列比对 (alignment)

- Aligning Millions of Short Sequence Reads
- One sequence is "embedded" in the other sequence (NGS reads, PCR primer, etc.)
  - ✦ "Local alignment" for long sequence
  - ✦ "Global alignment" for short sequence
- Aligners:
  - ✦ BWA, Bowtie2, STAR, HISAT2, …
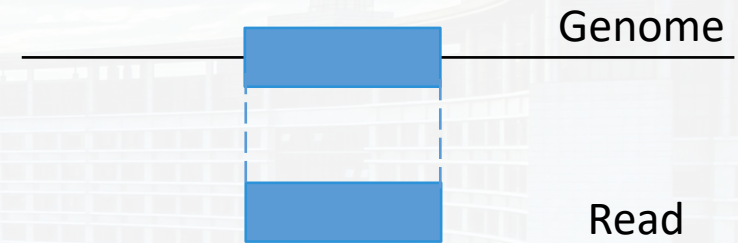  - ✦ minimap2, STARlong… (Nanopore, PacBio)

# Mapping: Input data

- Reference Genome
  - ✦ Nucleotide sequence (FastA)
  - ✦ Length: Hundreds of Mb per chromosome
  - ✦ ~3 Gb in total (for human genome)
- Reads
  - ✦ Nucleotide sequence with various qualities (FastQ)
    - error rate ranges from a few tenths of a percent to several percent
  - ✦ Length: ~100 bp per read
  - ✦ Hundreds of Gbs per run

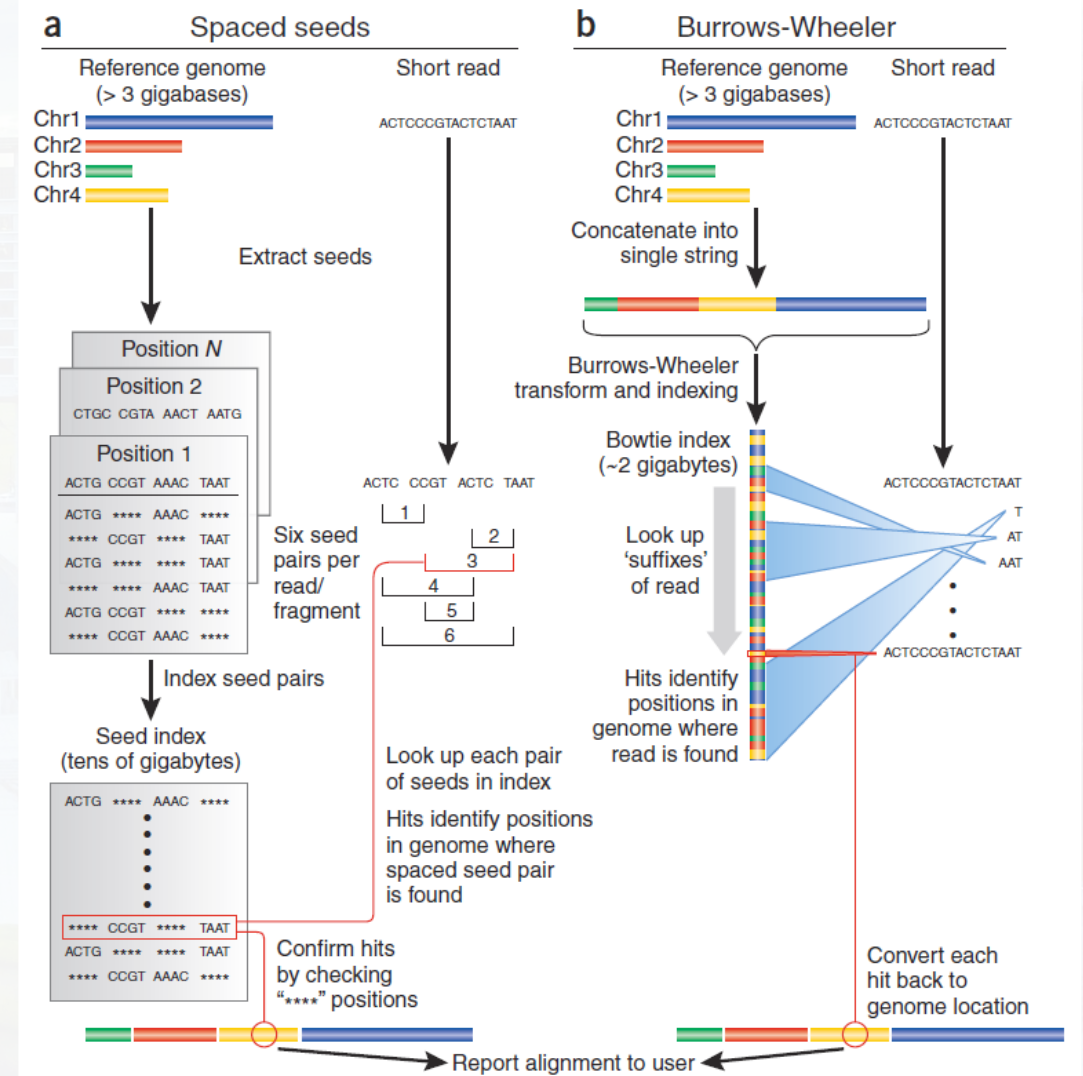# Mapping algorithms

- Burrows-Wheeler transform (BWT)
  - ✦数据压缩算法(bzip2)
  - ✦BWT-based tools: BWA, Bowtie2, SOAP2
  - ✦Fast, memory-efficient, Less sensitive

- Hashing (哈希)
  - ✦Hash-based tools: MAQ, Novoalign, Stampy
  - ✦most accurate overall results



Trapnell, C. & Salzberg, S. L. (2009) How to map billions of short reads onto genomes. Nat Biotech 27, 455–457.
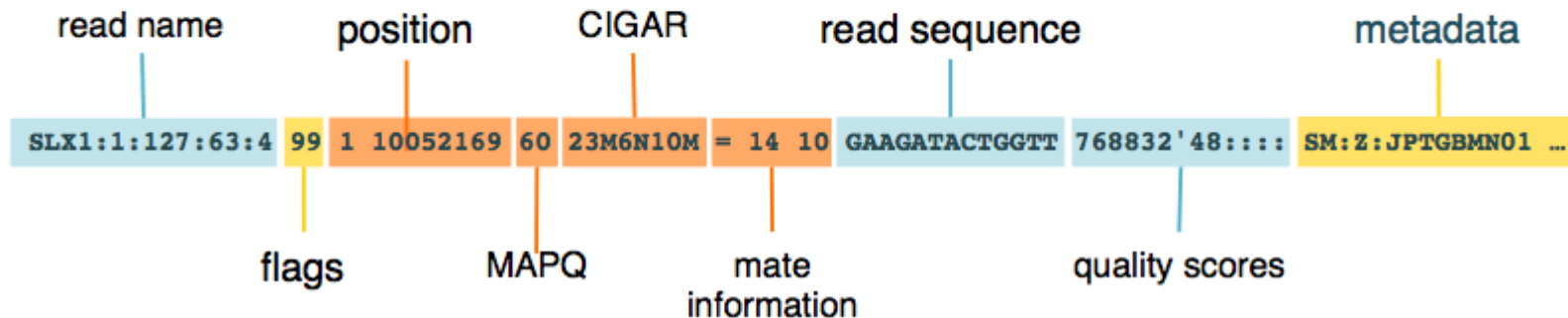
# 比对结果文件 — SAM

- SAM(Sequence Alignment/Map) format is a generic format for storing large nucleotide sequence alignments
  - ✦ 在SAM格式中，每一行表示 一个read的比对结果

```
@HD VN:1.6 SO:coordinate
@SQ SN:ref LN:45
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       * SA:Z:ref,29,-,6H5M,17,0;
r004    0 ref 16 30 6M14N5M    *  0   0 ATAGCTTCAGC       *
r003 2064 ref 29 17 6H5M       *  0   0 TAGGC             * SA:Z:ref,9,+,5S6M,30,1;
r001  147 ref 37 30 9M         =  7 -39 CAGCGGCAT         * NM:i:1
```

**HEADER** containing metadata (sequence dictionary, read group definitions etc)
**RECORDS** containing structured read information (1 line per read record)

| read name | | position | | CIGAR | | | | read sequence | | metadata |
|---|---|---|---|---|---|---|---|---|---|---|
| SLX1:1:127:63:4 | 99 | 1 10052169 | 60 | 23M6N10M | = 14 | 10 | | GAAGATACTGGTT | 768832'48:::: | SM:Z:JPTGBMNO1 ... |

flags    MAPQ    mate information    quality scores

Mapping quality (MAPQ) score is the probability that the read is incorrectly mapped, or more importantly, the probability that the read maps uniquely.
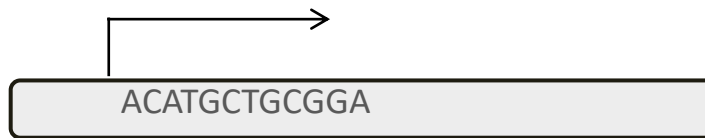
# Read Alignment
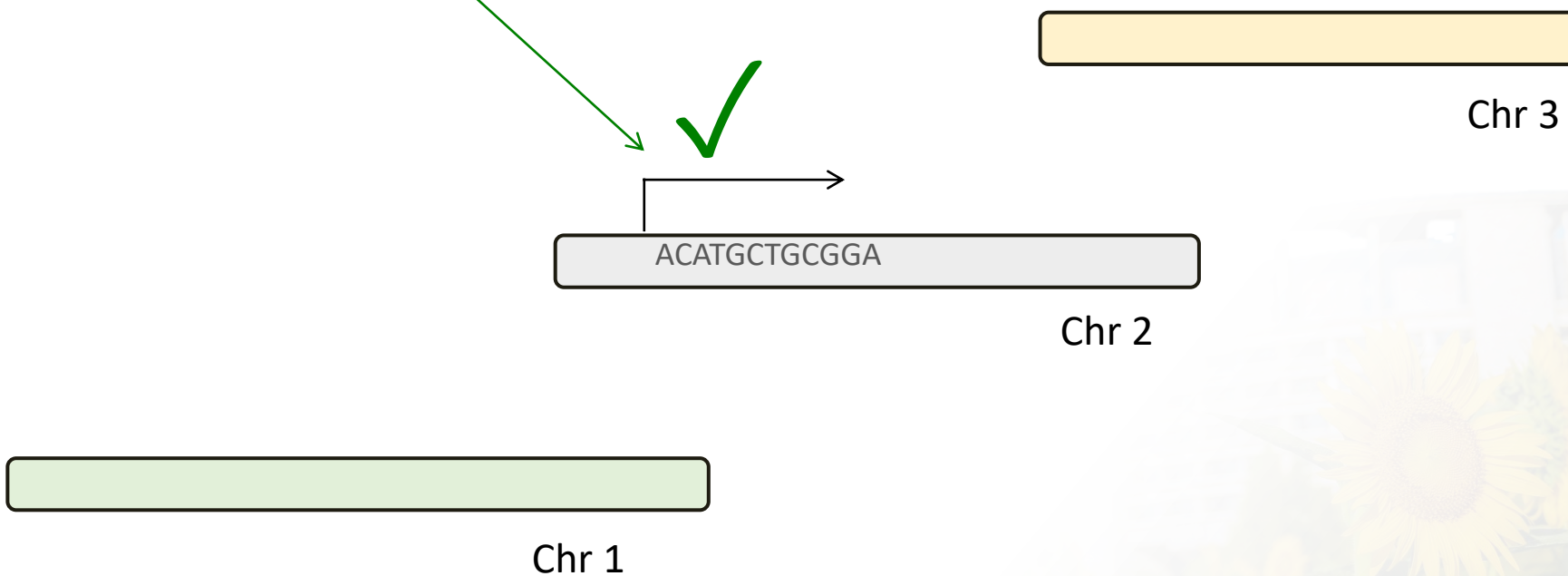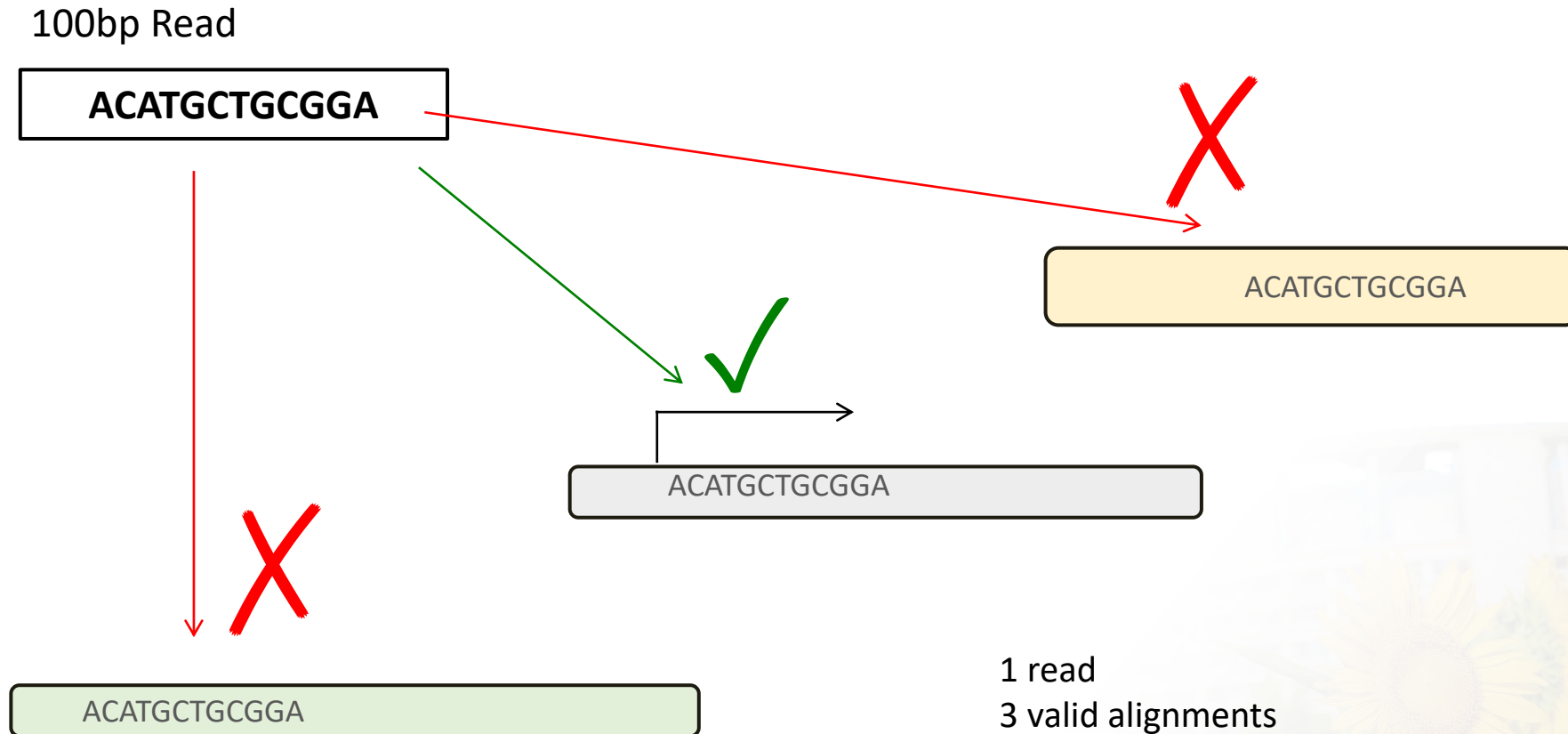
100bp Read

ACATGCTGCGGA

Reference sequence

Chr 3

ACATGCTGCGGA

Chr 2

Chr 1

Steve Munger, 2017

# The perfect read: 1 read = 1 unique alignment.



100bp Read

ACATGCTGCGGA

✓

ACATGCTGCGGA

Chr 3

Chr 2

Chr 1

Steve Munger, 2017

# Some reads will align equally well to multiple locations.  "Multi-mapped reads"

100bp Read

**ACATGCTGCGGA**

ACATGCTGCGGA

✓

ACATGCTGCGGA

ACATGCTGCGGA

1 read
3 valid alignments
Only 1 alignment is correct

- Ignore them?
- Weight them?

Steve Munger, 2017

# SAM文件工具 — SAMtools

- Tools to handle Bam/Sam files: SAMtools
  - ✦ $samtools view test.bam
  - ✦ $samtools view –h  test.bam | less  #show headers
  - ✦ $samtools flagstat ./data/SRR3096662_Aligned.sort.bam

- View alignment with samtools
  - ✦ $samtools index ./SRR3096662_Aligned.sort.bam
  - ✦ $samtools tview ./SRR3096662_Aligned.sort.bam --reference ./GRCh37.genome.fa
  
  #Need to make the index for the bam file



Aligned reads

# 比对结果可视化(Visualization):IGV

**Integrative Genome Viewer (IGV)：**

http://software.broadinstitute.org/software/igv/download

mapping

Reference Genome

...CCATAG       TAT  CGCCC      CGGA AATTT  CGGTATAC
...CCAT     CTATAT CG        TCGGA AATT   CGGTATAC
...CCAT  GGCTATAT CGC  CTATCGGAAA       GCGGTATA
...CCA  AGGCTATAT CGC CCTATCGGA       TTGCGGTA   C...
...CCA  AGGCTATAT    GCCCTATCG A        TTTGCGGT      C...
...CC    AGGCTATAT    GCCCTATCG  AAATTTGC      ATAC...
...CC  TAGGCTATA   CGCCCTA       AAATTTGC  GTATAC...
...CCATAGGCTATATGCGCCCTATCGGCAATTTGCGGTATAC...

Genetic variants

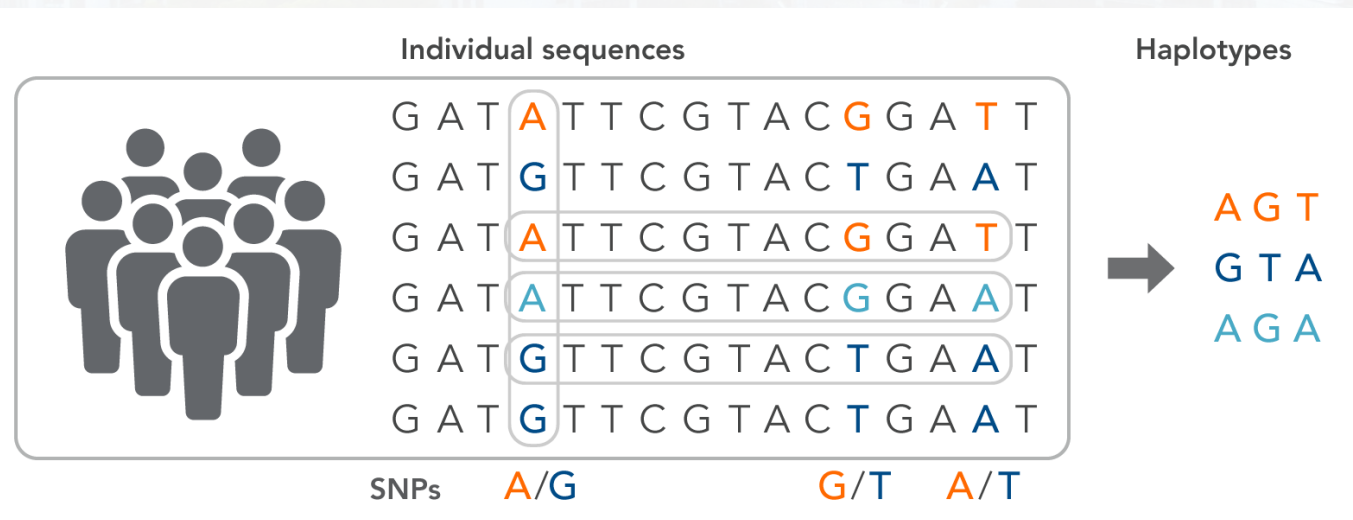# SNP calling VS. Genotyping

- SNP calling: identifies variable sites (variants).
- Genotyping: determines the genotype for each individual at each site.
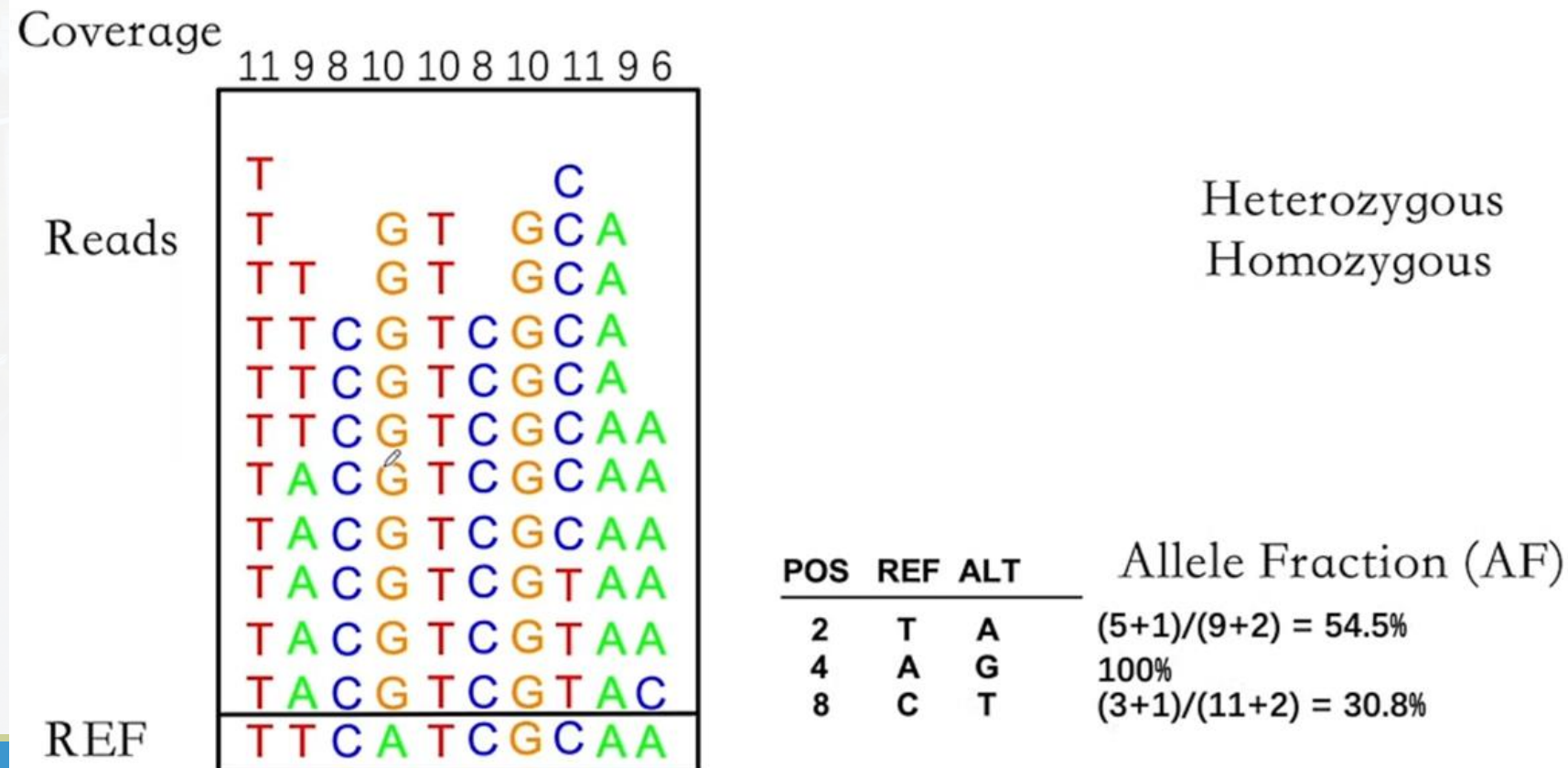  - ✦ The number of alleles (等位基因) or ploidy (染色体倍性) is decided and fixed.

浙江工商大学
ZHEJIANG GONGSHANG UNIVERSITY

- The sequenced DNA fragments ("reads") are aligned to the reference sequence and the base at each position is determined, counted and then run through a number of statistical tests to determine whether the site is different from the reference sequence.



Coverage
11 9 8 10 10 8 10 11 9 6

Reads

Heterozygous
Homozygous

| POS | REF | ALT | Allele Fraction (AF) |
|-----|-----|-----|----------------------|
| 2 | T | A | (5+1)/(9+2) = 54.5% |
| 4 | A | G | 100% |
| 8 | C | T | (3+1)/(11+2) = 30.8% |

Confidence?

# Algorithms for genotype and SNP calling

- SNP calling can be done using likelihood ratio tests or Bayesian procedures, …
- 遗传变异的统计学方法可以参考文献:

Review Article | Published: 18 May 2011

## Genotype and SNP calling from next-generation sequencing data

Rasmus Nielsen ✉, Joshua S. Paul, Anders Albrechtsen & Yun S. Song ✉

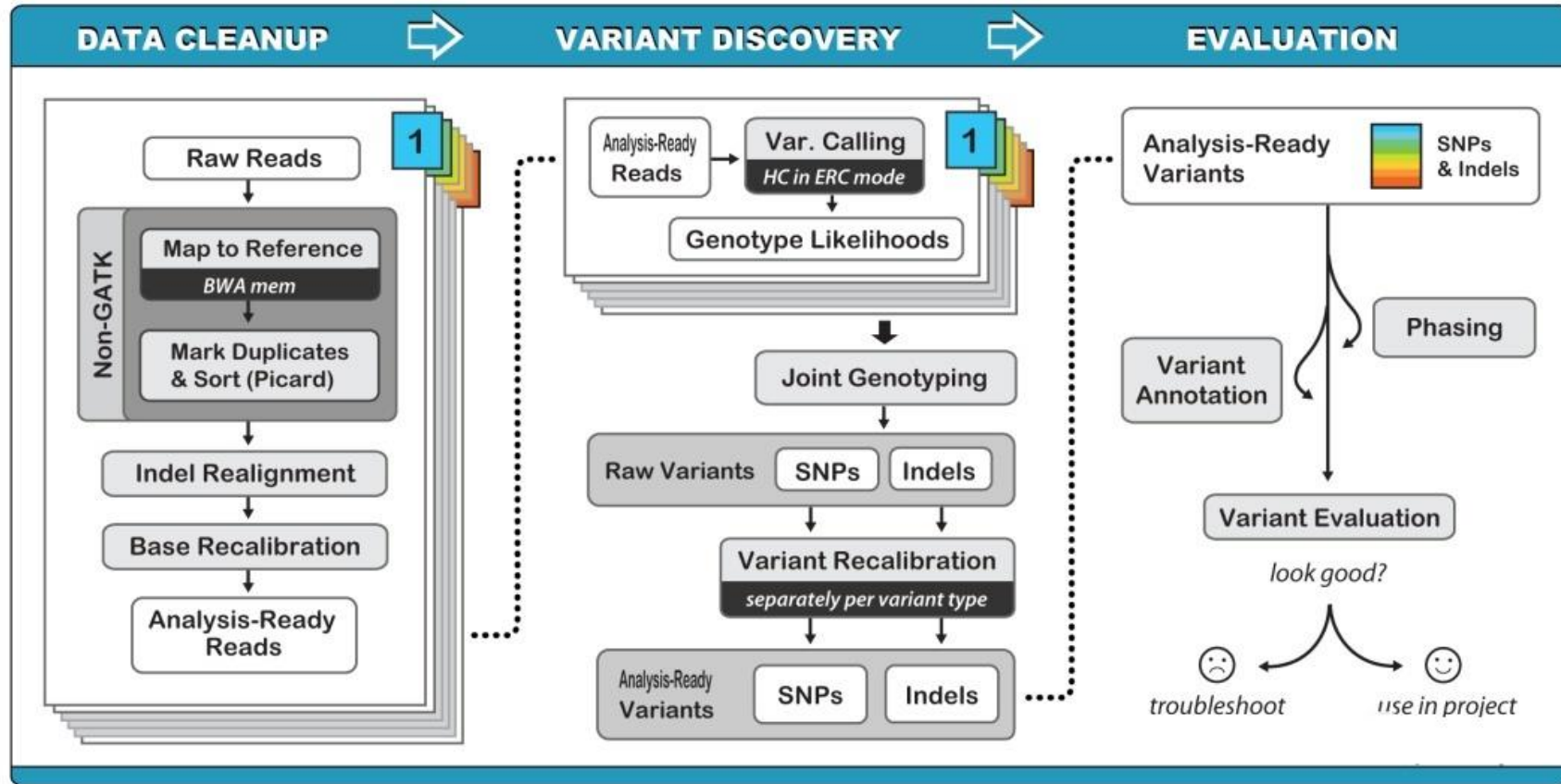*Nature Reviews Genetics* **12**, 443–451 (2011) | Cite this article

**42k** Accesses | **907** Citations | **34** Altmetric | Metrics

## Key Points

- Converting next-generation sequencing (NGS) image files into a set of called SNPs involves a number of steps including image analysis, alignment and assembly, SNP calling and genotype calling.

- Genotype probabilities for a single individual can be calculated from alignments using recalibrated quality scores.

Nielsen, R., Paul, J., Albrechtsen, A. et al. Genotype and SNP calling from next-generation sequencing data. Nat Rev Genet 12, 443–451 (2011).

# 常用的变异鉴定工具:
# GATK (Genome Analysis ToolKit)



best practices:
https://gatk.broadinstitute.org/hc/en-us/articles/360036194592-Getting-started-with-GATK4

- VCF(Variant Call Format) 是存储遗传变异类型，如SNPs, Indels和SVs的标准文件格式
  - ✦ 在VCF格式中，每一行表示 一个遗传变异的结果

**Example VCF file**

```
##fileformat=VCFv4.2
##FORMAT=<ID=GT,Number=1,Type=Integer,Description="Genotype">
##FORMAT=<ID=GP,Number=G,Type=Float,Description="Genotype Probabilities">
##FORMAT=<ID=PL,Number=G,Type=Float,Description="Phred-scaled Genotype Likelihoods">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT SAMP001 SAMP002
20 1291018 rs11449 G A . PASS . GT 0/0 0/1
20 2300608 rs84825 C T . PASS . GT:GP 0/1:. 0/1:0.03,0.97,0
20 2301308 rs84823 T G . PASS . GT:PL ./.:. 1/1:10,5,0
```

http://samtools.github.io/hts-specs/

# 基因组变异的检测结果-VCF



以#开头的注释部分

主体部分

1：参考序列名称　　4：参考序列碱基　　7：位点是否需要过滤

2：变异所在位置　　5：目的碱基　　　　8：变异相关信息

3：变异点ID　　　　6：变异质量值　　　9：变异格式

# SARS-CoV-2基因组重测序分析

- WuHan-Hu-1

- UK, South African, Brazilian, and India variants

- The World Health Organization (WHO) has designated SARS-CoV-2 variants Alpha(B.1.1.7), Beta(B.1.351), Gamma(P.1), and Delta(B.1.617.2) as Variants of Concern (VOC), among which Delta variant with remarkable transmission and immune escape ability has attracted great attentions.

基因组的保守区与易变区



Total number of amino acid substitutions found in **4,400 coronavirus genomes** from Dec. to April

Longer lines may → show places where the genome is more tolerant of mutations.

Gaps may show critical spots in the genome that cannot tolerate mutations.

哪个区适合适作为抗病毒药物或疫苗的靶标？

- Mutations in the genome produce a fingerprint that can be used to infer ancestral relationships (phylogeny).

- Zoonosis, jumped from animals to humans



Transmission Cycle of SARS CoV 2. Contributed by Rohan Bir Singh, MD; Made with Biorender.com

## 1.1 软件工具的准备

(1)操作系统: Linux/macOS

(2)测序序列质量控制软件: FastQC、Trimmomatic

(3)序列比对工具: BWA，包括BWA index、BWA mem、BWA aln、BWA sample

(4)序列分析工具: Samtools，包括Samtools view、Samtools sort、Samtools mpileup、 bcftools call、bcftools consensus

## 1.2 数据的准备

### (1)参考序列

- SARS-CoV-2 Genome (NC_045512.2)

### (2)测序数据(Illumina PE)

- NCBI accession number: SRR11140750

- A clinical swab obtained from a patient in Madison, WI, USA

- Data released on 21$^{st}$ Feb 2020.

SRX7777160: **SARS-CoV-2 swab_illumina**

1 ILLUMINA (Illumina MiSeq) run: 17,657 spots, 7.7M bases, 3.6Mb downloads

**Design:** SISPA Nextera XT

**Submitted by:** University of Wisconsin - Madison

**Study:** SARS-CoV-2 parallel sequencing by Illumina and Oxford Nanopore Technologies

PRJNA607948 • SRP250294 • All experiments • All runs

## 2.1 FastQC 用于测序数据概况的分析

$mkdir fastqc_out
$fastqc SRR11140750.fastq -o ./fastqc_out/

-o --outdir   FastQC生成的报告文件的储存路径；

　--extract　 默认情况下生成的报告打包成1个压缩文件，设置该参数是无需打包；

-t –threads　 选择程序运行的线程数，每个线程占用250MB内存，越多程序运行越快

-c --contaminants　污染物选项，输入的是一个文件，文件格式是 Name [Tab] Sequence，里面是可能的污染序列，如果有这个选项，FastQC会在计算时评估污染的情况，并在统计的时候进行分析；

-a –adapters　　也是输入一个文件，文件格式是 Name [Tab] Sequence，储存的是测序的adpater序列信息；如果不输入，当前版本的 FastQC 会按照通用引物来评估序列是否有 adapter的残留。

## 2.2 除去接头和低质量序列: Trimmomatic

$java -jar Trimmomatic-0.39/trimmomatic-0.39.jar SE -phred33 SRR11140750.fastq SRR11140750.out.fq ILLUMINACLIP:Trimmomatic-0.39/adapters/TruSeq3-SE.fa:2:30:10 SLIDINGWINDOW:5:20 LEADING:5 TRAILING:5 MINLEN:50

-ILLUMINACLIP: 设置切除接头序列的参数

-SLIDINGWINDOW: 设置滑动窗口长度的参数，

-LEADING: 设置是否切除read开头碱基的质量阈值；

-TRAILING: 设置是否切除read末尾碱基的质量阈值；

-MINLEN: 设置read被切除后至少需要保留的长度；如果低于该长度，序列将被丢弃。

\#参考序列NC_045512.2文件重命名为covid19.fasta

$cp NC_045512.2.fasta covid19.fasta

\#Index the reference genome for use with BWA

$bwa index covid19.fasta

\#Align the Illumina reads

$bwa mem covid19.fasta SRR11140750.fastq.gz > SRR11140750.sam

\#Coordinate sort SAM file, and output to BAM

$samtools sort -o SRR11140750.bam SRR11140750.sam

Can you figure out how to do the bwa mem and samtools sort commands
in a pipeline so as to avoid writing the large intermediary SAM file?

**#Generate index of the genome file**

$samtools faidx covid19.fasta

**#Index the BAM file**

$samtools index SRR11140750.bam

**# View alignment with samtools**

$samtools tview ./SRR11140750.bam --reference ./covid19.fasta

#Press "q" to quit view

最终生成四个文件：

- COVID19.fasta
- COVID19.fasta.fai
- SRR11140750.bam
- SRR11140750.bam.bai

# 4.BAM file visualization with IGV

- Run local IGV, or visit IGV-web (https://igv.org/app/).

- Load the genome from a "Local File ..." by selecting both the COVID-19.fasta and COVID-19.fasta.fai files.

- Load a "Track" from a "Local File ..." by selecting both the SRR11140750.bam and SRR11140750.bam.bai files.



The stacked grey arrows represent the reads aligned to the SARS-CoV-2 reference genome. What do the coloured vertical bars within the reads indicate?

- 使用本课提供的SARS-CoV-2测序数据(SRR11140750.fastq)，利用 BWA、Samtools、IGV等工具进行reads mapping与可视化等。

谢 谢 观 看