

# Variant Calling

**BCB 5200 Introduction Bioinformatics I**

Fall 2017

Tae-Hyuk (Ted) Ahn

Department of Computer Science  
Program of Bioinformatics and Computational Biology  
Saint Louis University



**SAINT LOUIS  
UNIVERSITY™**

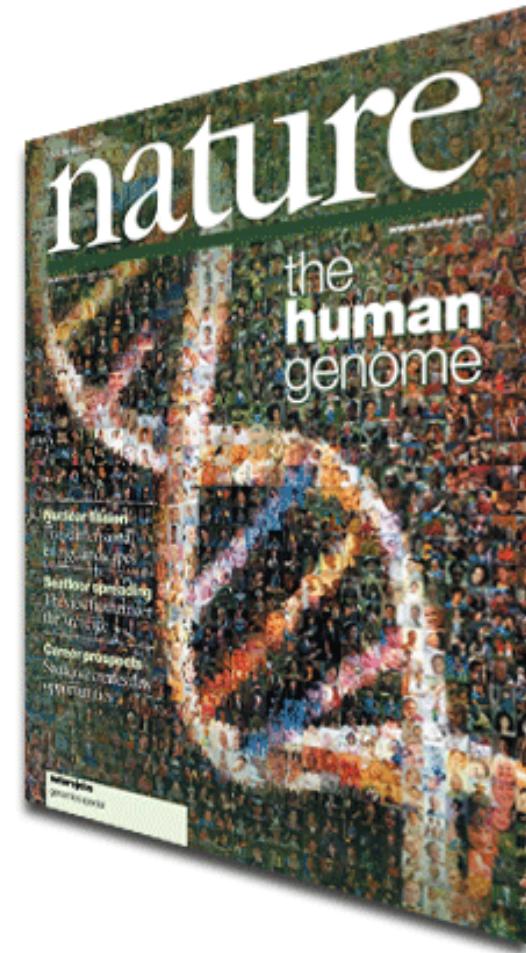
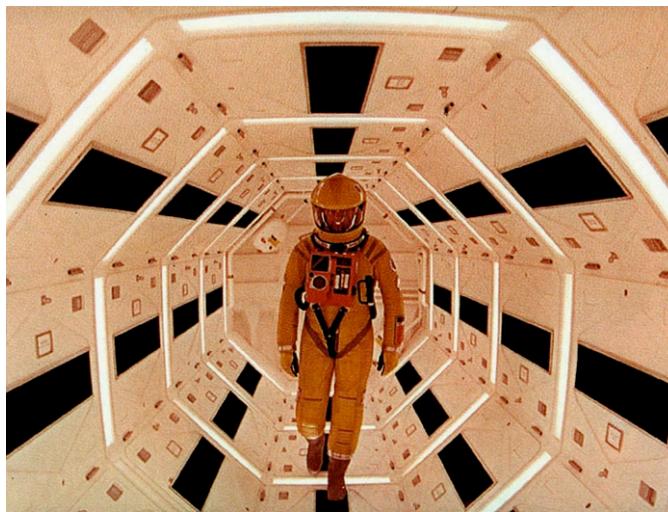
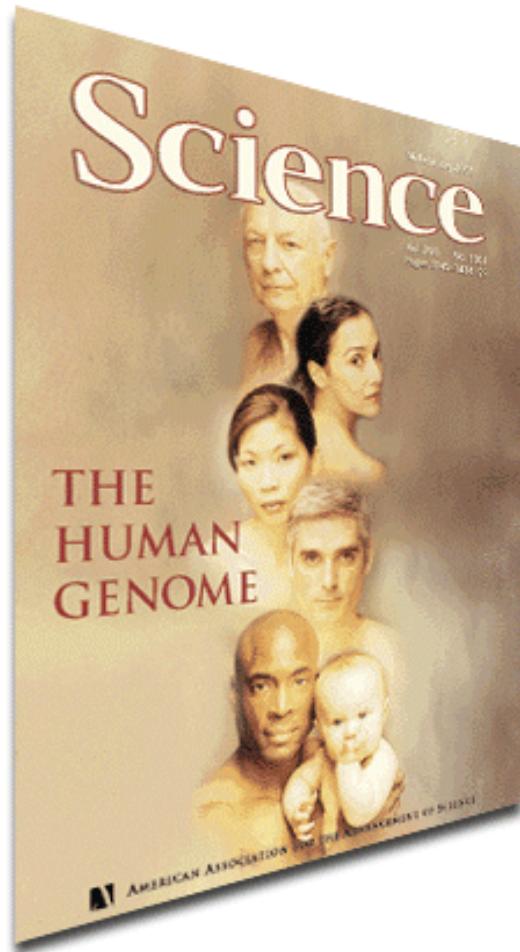
— EST. 1818 —

# The Human Genome

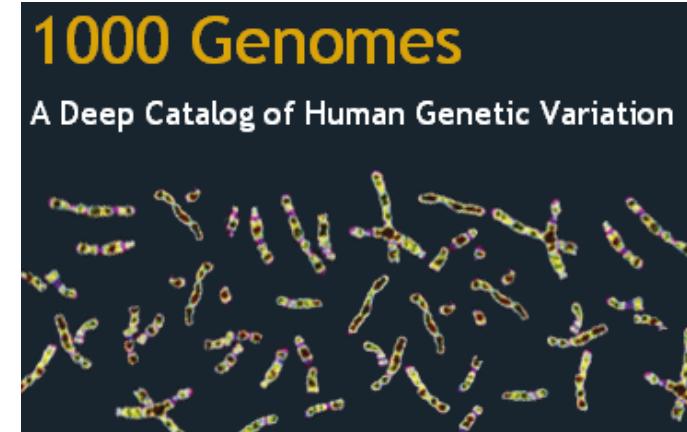
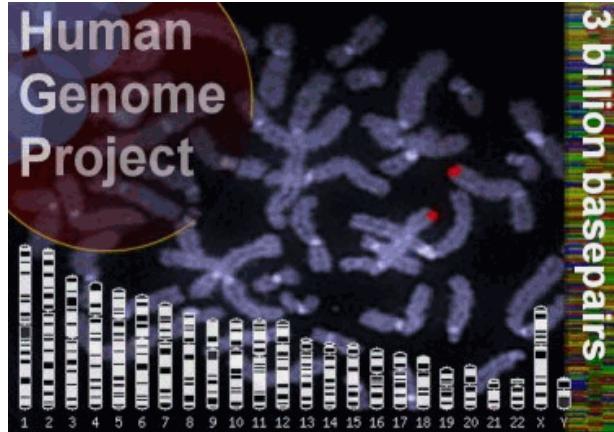
- 3 billions base pairs (**ATGC**)
- 20,000 protein-coding genes
- 99.6% inter-individual identity (yet 4 millions differences)
- 99% identical to chimpanzee genome (yet 6% different genes)



# The Human Genome



# Exploring the Human Genome



2002

2008



# International HapMap Project

- The International HapMap Project (launched in 2002) was an organization that aimed to develop a haplotype map (HapMap) of the human genome, to describe the common patterns of human genetic variation.
- HapMap is used to find genetic variants affecting health, disease and responses to drugs and environmental factors.
- Canada, China, Japan, Nigeria, the United Kingdom, and the United States.



# 1000 Genome Project

- Whole genome sequencing and complete description of human genetic diversity in >1000 individuals from multiple world populations
- <http://www.1000genomes.org>

The 1000 Genomes Project is the first project to sequence the genomes of a large number of people, to provide a comprehensive resource on human genetic variation.



# What is Cancer Genomics?

- In cancer cells, small changes in the genetic letters can change what a genomic word or sentence means.
- A changed letter can cause the cell to make a protein that doesn't allow the cell to work as it should.
- These proteins can make cells grow quickly and cause damage to neighboring cells.
- By studying the cancer genome, scientists can discover what letter changes are causing a cell to become a cancer.
- The genome of a cancer cell can also be used to tell one type of cancer from another.

# Short Video

## How to sequence the human genome - Mark J. Kiel

<http://ed.ted.com/lessons/how-to-sequence-the-human-genome-mark-j-kiel>

# The Cancer Genome Atlas (TCGA)

- The Cancer Genome Atlas (TCGA) is a project, begun in 2005, to catalogue genetic mutations responsible for cancer, using genome sequencing and bioinformatics.
- TCGA applies high-throughput genome analysis techniques to improve our ability to diagnose, treat, and prevent cancer through a better understanding of the genetic basis of this disease.
- <http://cancergenome.nih.gov/>

The screenshot shows the homepage of The Cancer Genome Atlas (TCGA). At the top left is the NIH logo followed by the text "THE CANCER GENOME ATLAS" and "National Cancer Institute National Human Genome Research Institute". At the top right are links for "Launch Data Portal", "Contact Us", and "For the Media". Below the header is a search bar with a magnifying glass icon and a "Search" button. A navigation bar at the bottom includes links for "Home", "About Cancer Genomics", "Cancers Selected for Study", "Research Highlights", "Publications", "News and Events", and "About TCGA".

# Types of Genetic Variation

Cancer is driven by genomic alterations like:

- Single Nucleotide Aberrations
  - Single Nucleotide Polymorphisms (**SNPs**) - mutations shared amongst a population
  - Single Nucleotide Variations (**SNVs**) - private mutations
- Short Insertions or Deletions (**indels**)
- Copy Number Variations (**CNVs**)
- Larger Structural Variations (**SVs**)

# SNPs vs. SNVs

Both are aberrations at a single nucleotide

- **SNP**

- Aberration expected at the position for any member in the species (well-characterized)
- Occur in population at some frequency so expected at a given locus
- Validated in population
- Catalogued in dbSNP (<http://www.ncbi.nlm.nih.gov/snp>)

- **SNV**

- Aberration seen in only one individual (not well characterized)
- Occur at low frequency so not common
- Not validated in population

Really a matter of frequency of occurrence

# SNVs of interest

- **Non-synonymous** mutations
  - Result in amino acid change
  - Impact protein sequence
  - Missense, nonsense, stop gained/lost mutations
- **Somatic** mutations in cancer
  - Tumor-specific mutations

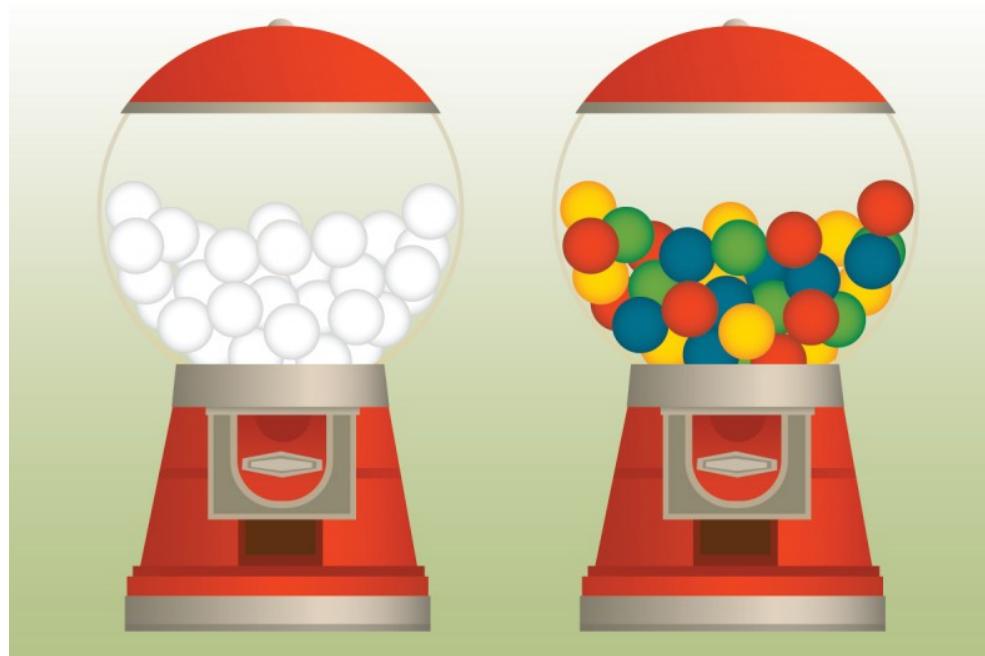
# Challenges of accurate somatic variant calling

Not as simple as identifying sites with a variant allele in the tumor not present in the normal

- Artifacts from PCR amplification or targeted (exome) capture
- Machine sequencing errors
- Incorrect local alignment of reads
- Tumor heterogeneity
- Tumor-normal cross-contamination

# What Is Tumor Heterogeneity?

- Until recently, the cells within a tumor were thought to be similar to one another at any given stage of the cancer, like the white gumballs in the machine on the left.
- But scientists are finding that a person's **tumor cells can be highly diverse**. This could mean that a biopsy may capture only a fraction of tumor cells that are not representative of the whole tumor, like a ball ejected at random from the machine on the right

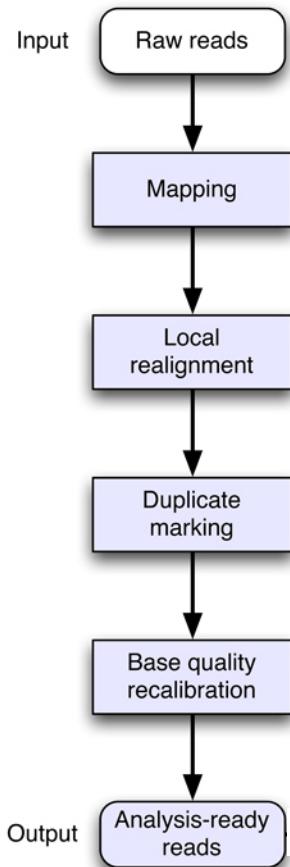


<https://www.mskcc.org/blog/what-tumor-heterogeneity>

# A framework for variation discovery

## Phase 1: NGS data processing

— Typically by lane —



## Phase 1: Mapping

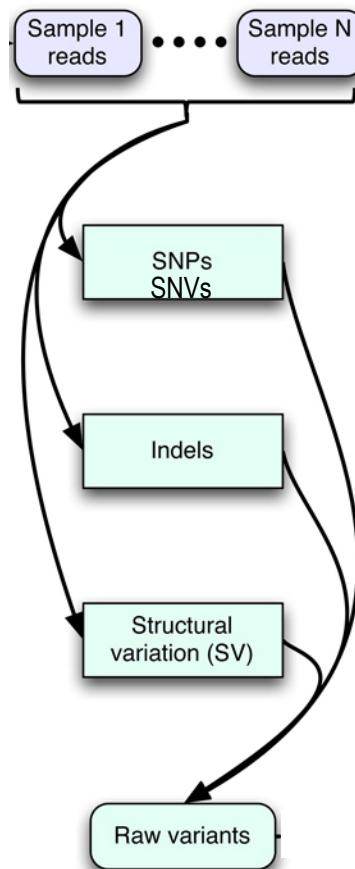
- Place reads with an initial alignment on the reference genome using mapping algorithms
- Refine initial alignments
  - local realignment around indels
  - molecular duplicates are eliminated
- Generate the technology-independent SAM/BAM alignment map format

Accurate mapping crucial for variation discovery

DePristo, M.A. et al. A framework for variation discovery and genotyping using next-generation DNA sequencing data. Nat Genet. 43(5):491-8. PMID: 21478889 (2011).

# A framework for variation discovery

## Phase 2: Variant discovery and genotyping

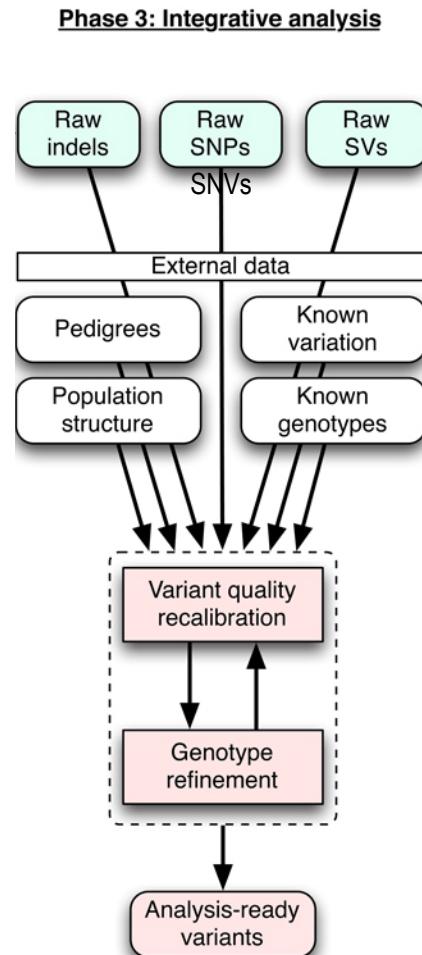


## Phase 2: Discovery of raw variants

- Analysis-ready SAM/BAM files are analyzed to discover all sites with statistical evidence for an alternate allele present among the samples
- SNPs, SNVs, short indels, and SVs

DePristo, M.A. et al. A framework for variation discovery and genotyping using next-generation DNA sequencing data. Nat Genet. 43(5):491-8. PMID: 21478889 (2011).

# A framework for variation discovery

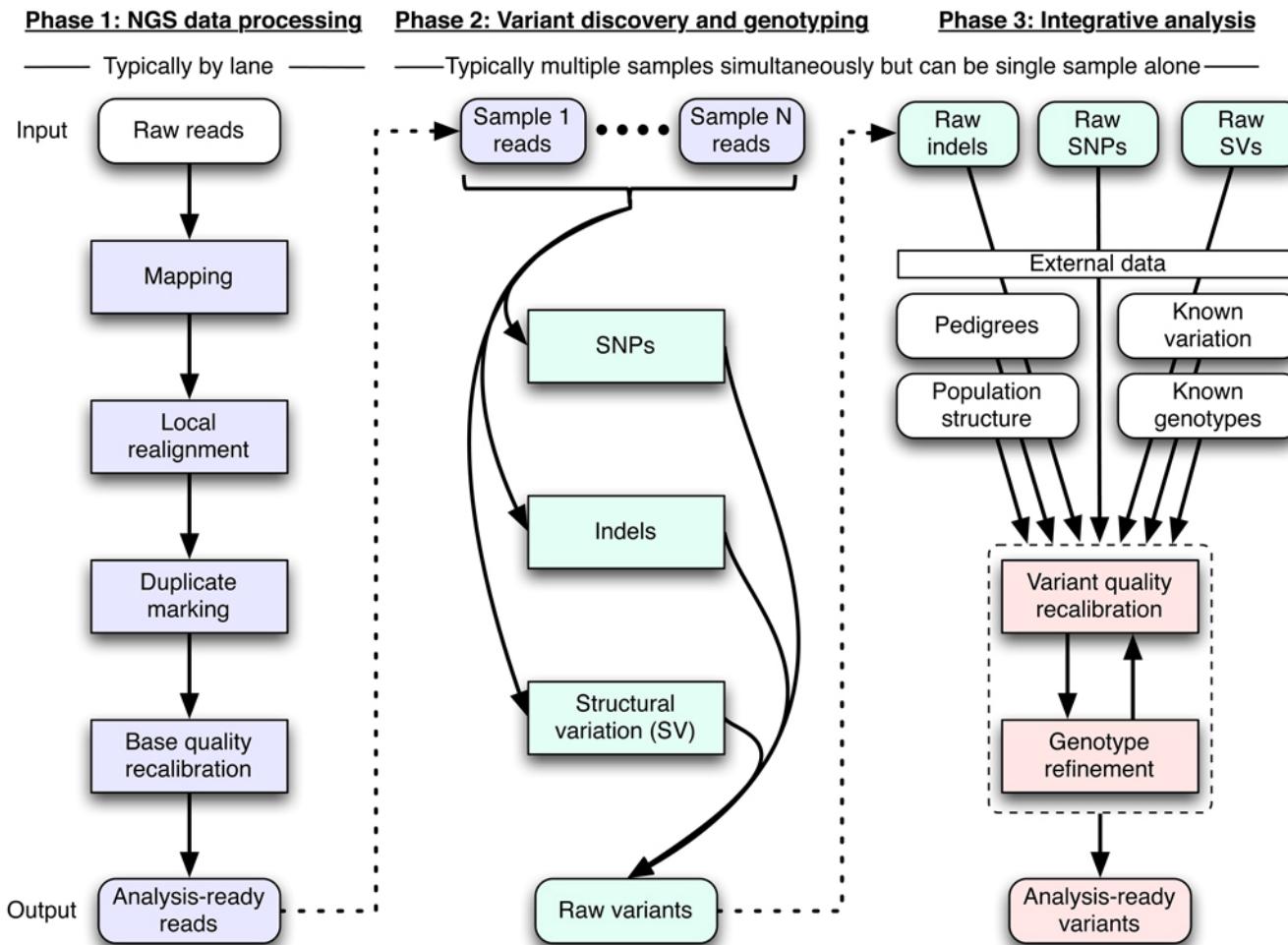


## Phase 3: Discovery of analysis-ready variants

- technical covariates, known sites of variation, genotypes for individuals, linkage disequilibrium, and family and population structure are integrated with the raw variant calls from Phase 2 to separate true polymorphic sites from machine artifacts
- at these sites high-quality genotypes are determined for all samples

DePristo, M.A. et al. A framework for variation discovery and genotyping using next-generation DNA sequencing data. Nat Genet. 43(5):491-8. PMID: 21478889 (2011).

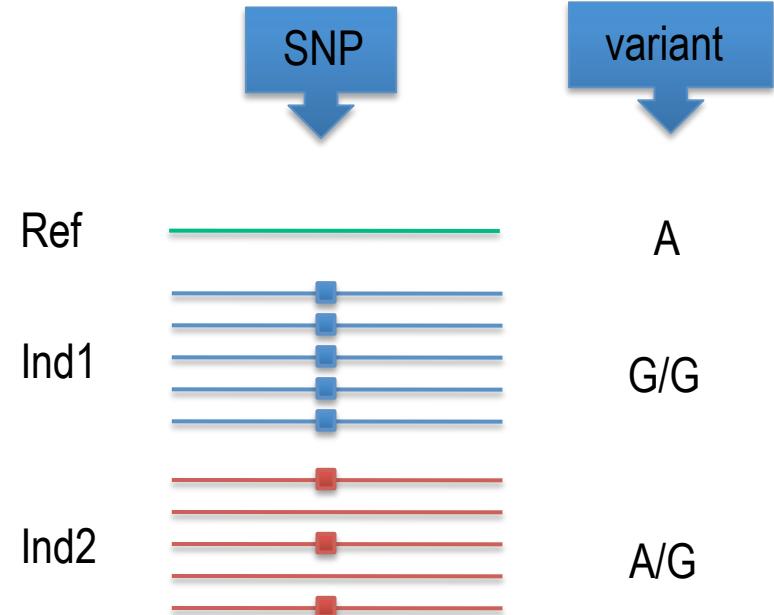
# A framework for variation discovery



DePristo, M.A. et al. A framework for variation discovery and genotyping using next-generation DNA sequencing data. Nat Genet. 43(5):491-8. PMID: 21478889 (2011).

# Variant calling methods

- > 15 different algorithms
- Three categories
  - Allele counting
  - Probabilistic methods, e.g. Bayesian model
    - to quantify statistical uncertainty
    - Assign priors based on observed allele frequency of multiple samples
  - Heuristic approach
    - Based on thresholds for read depth, base quality, variant allele frequency, statistical significance



Nielsen R, Paul JS, Albrechtsen A, Song YS. Genotype and SNP calling from next-generation sequencing data. Nat Rev Genet. 2011 Jun; 12(6):443-51. PMID: 21587300.

<http://seqanswers.com/wiki/Software/list>

# Sequencing

Library construction and sequencing

Sequencing quality control  
FASTQ files

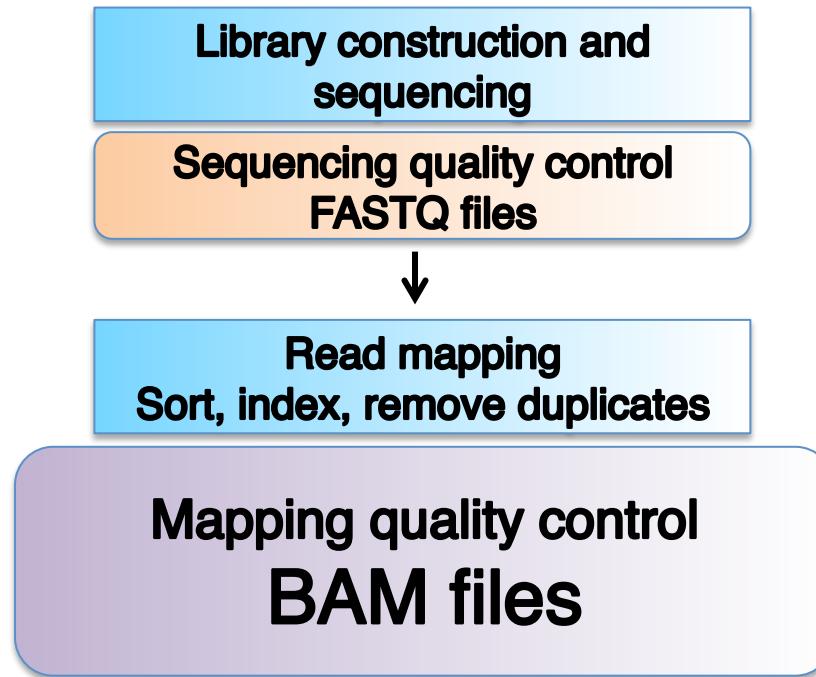
@G:1:1:11:1079#0/1

TGATTGATTCCATTCCATTCCATTCCATTCCATTGCAATCCCTCCAATCCATTCCATTCCATTCCATTTC

+G:1:1:11:1079#0/1

`Xa^YO\\_\^a\\_`\_\_`a\_\_^a^a^\_a``^\_\\`\\]``[XUGXXXXXWUTWWVWUSTXXPUWYYRVWYYYZYXYWZ

# Mapping



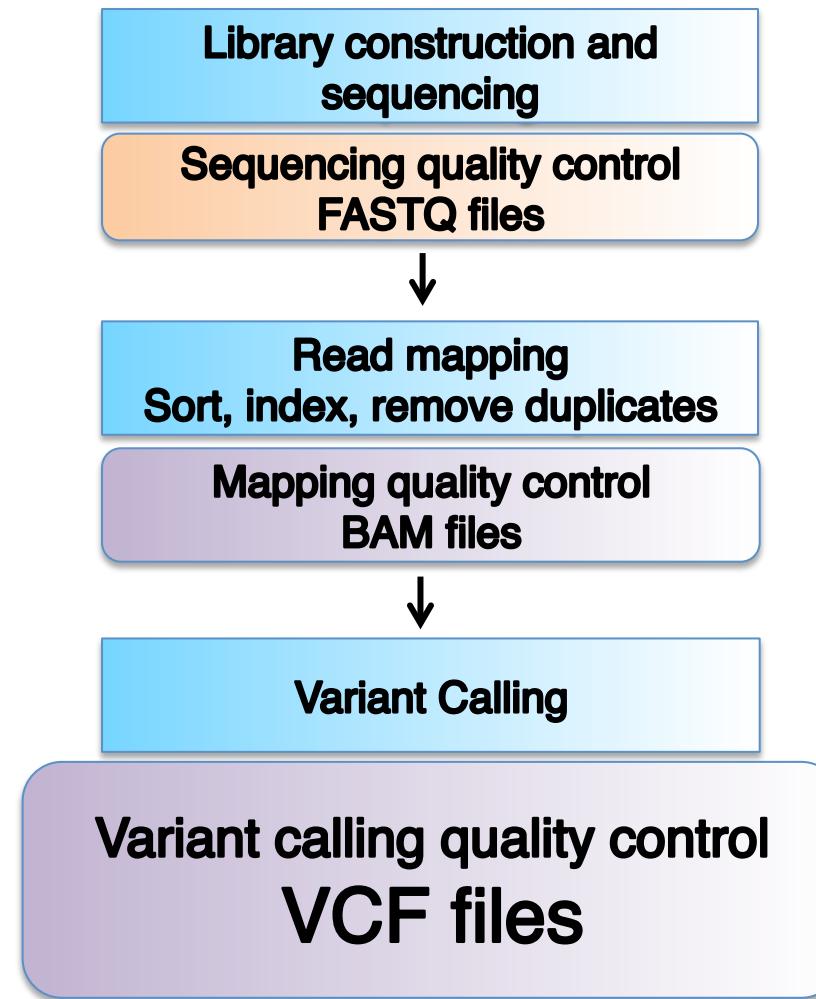
# Mapping

711      721      731      741      751      761      771      781      791

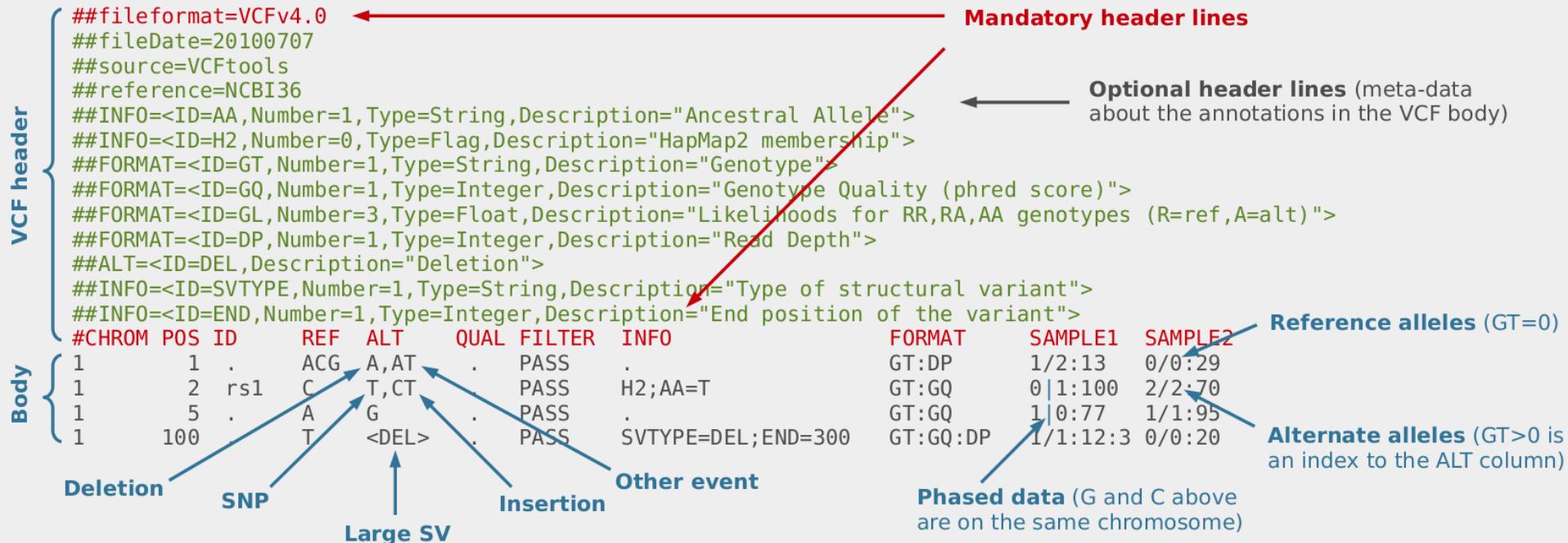
GTGGCGAGAAAATGTCGATGCCATTATGCCGGCGTGTAGAACGCGTGGTACAACGTTACCGTTATCGATCCGGTCGAAAAACTGCTGGCAGT

.....G.....  
.....g.....c.....  
,,c,,t,,g,,.....  
,,t,,g,,.....  
,,g,,.....  
,,tg,,.....,t,,  
,,c,,g,,g,,.....  
c,,g,,a,,.....,t,,  
,,gc,,.....  
,,g,g,,.....  
,,g,,g,,.....  
,,g,,g,,.....,a,c,,  
,,g,,g,,.....,g,,.....  
,,g,,g,,.....c,,  
,,g,,.....  
,,g,,.....  
,,g,,g,,.....,t,,  
,,g,,g,,.....,t,,  
,,g,,g,,.....,c,,  
,,g,,g,,.....,a,,g,,  
,,c,,g,,.....,a,,t,,  
,,tg,,g,,.....,g,,.....,g,,.....,a,,

# Variant Calling



# Variant Call Format (VCF)



# Tools for Variant Calling

- Alignmer
  - Bowtie2 (<http://bowtie-bio.sourceforge.net/bowtie2/>)
  - BWA (<http://bio-bwa.sourceforge.net/>)
- Tools for variant calling
  - SAMtools & BCFtools (<http://www.htslib.org/download/>)
  - VCFTools (<https://vcftools.github.io/index.html>)
  - GATK (<https://www.broadinstitute.org/gatk/>)
- Browser
  - SAMtools
  - IGV (<https://www.broadinstitute.org/igv/>)
  - UCSC Genome Browser (<https://genome.ucsc.edu/>)

# Variant Calling Lab

1 Generate a small set of simulated reads for *E. coli*.

wgsim

2 Align the reads to the reference *E. coli* genome.

bowtie2

3 Convert the aligned reads from the SAM file format to BAM.

4 Sort and index the BAM file.

samtools

5 Identify genomic variants.

6 Visualize the reads and genomic variants.

# Generating Simulated Reads

- To download wgsim (<https://github.com/lh3/wgsim>)
  - You can download zip file
  - You can download as

```
$ git clone https://github.com/lh3/wgsim.git
```
  - You don't need to install if the wgsim is on the system. Check it.
- How can you add the software path into your environment to run the program easily?
  - Open your ~/.bashrc profile and add the software path

```
$ pwd  
/faculty/ahnt/Courses/BCB5250/software/wgsim  
$ echo 'PATH=/faculty/ahnt/Courses/BCB5250/software/wgsim:$PATH' >> ~/.bashrc  
$ . ~/.bashrc
```

# Generating Simulated Reads

- First, we need a small set of sample read data.
- The command line usage for wgsim is:
- `wgsim [options] <in.ref.fa> <out.read1.fq> <out.read2.fq>`
- By default, wgsim reads in a reference genome in the **FASTA** format, and generates simulated **paired-end** reads in the **FASTQ** format.
- E. coli genome (NC\_008253.fna) is at /tmp directory. Copy the file into your Lab02\_VariantCalling directory.

```
$ wgsim -N 1000000 NC_008253.fna sim_reads_1M.fq
```

# Aligning Reads using Bowtie2

- Download Bowtie2 (before download, check it is already installed)

```
$ wget https://downloads.sourceforge.net/project/bowtie-bio/  
bowtie2/2.3.0/bowtie2-2.3.0-linux-x86_64.zip
```

- bowtie2-build to index your genome

```
$ bowtie2-build -f NC_008253.fna NC_008253.fna
```

- Run bowtie2

```
$ bowtie2 -x NC_008253.fna -U sim_reads_1M.fq -S  
sim_reads_1M_aligned.sam
```

# Understanding SAM Format

- As SAMtools is primarily concerned with manipulating SAM files, it is useful to take a moment to examine the sample SAM file generated by bowtie, and to dive into the details of the SAM file format itself. The first six lines from the bowtie SAM file are extracted below:

Each row describes a single alignment of a raw read against the reference genome. Each alignment has 11 mandatory fields, followed by any number of optional fields.

# Understanding SAM Format

| Field Name | Description   | Example from the <i>e_coli</i> SAM file   |
|------------|---|---|
| QNAME      | Unique identifier of the read; derived from the original FASTQ file.  | gi 110640213 ref NC_008253.1 _418_952_1:0:0_1:0:0_0/1   |
| FLAG       | A single integer value (e.g. 16), which encodes multiple elements of meta-data regarding a read and its alignment. Elements include: whether the read is one part of a paired-end read, whether the read aligns to the genome, and whether the read aligns to the forward or reverse strand of the genome. A <a href="#">useful online utility</a> decodes a single SAM flag value into plain English.  | 16: indicates that the read maps to the reverse strand of the genome.   |
| RNAME      | Reference genome identifier. For organisms with multiple chromosomes, the RNAME is usually the chromosome number; for example, in human, an RNAME of "chr3" indicates that the read aligns to chromosome 3. For organisms with a single chromosome, the RNAME refers to the unique identifier associated with the full genome; for example, in <i>E. coli</i> , the RNAME is set to: gi 110640213 ref NC_008253.1 .   | gi 110640213 ref NC_008253.1  |
| POS        | Left-most position within the reference genome where the alignment occurs.  | 418   |
| MAPQ       | Quality of the genome mapping. The MAPQ field uses a <b>Phred-scaled probability</b> value to indicate the probability that the mapping to the genome is incorrect. Higher values indicate increased confidence in the mapping.   | 42: indicates low probability (6.31e-05) that the mapping is incorrect.   |
| CIGAR      | A compact string that (partially) summarizes the alignment of the raw sequence read to the reference genome. Three core abbreviations are used: M for alignment match; I for insertion; and D for Deletion. For example, a CIGAR string of 5M2I63M indicates that the first 5 base pairs of the read align to the reference, followed by 2 base pairs, which are unique to the read, and not in the reference genome, followed by an additional 63 base pairs of alignment. Note that the CIGAR string is a partial representation of the alignment because M indicates an alignment match, but this could indicate an exact sequence match or a mismatch [2]. If you would like to determine match v. mismatch, you can consult the optional MD field, detailed below. | 70M: 70 base pairs within the read match the reference genome.  |
| RNEXT      | Reference genome identifier where the mate pair aligns. Only applicable when processing paired-end sequencing data. For example, an alignment with RNAME of "chr3" and RNEXT of "chr4" indicates that the mate and its pair span two chromosomes, indicating a possible structural rearrangement. A value of * indicates that information is not available.   | the <i>E. coli</i> simulated data is single read sequencing data, and all RNEXT values are therefore set to * (no information available). |
| PNEXT      | Position with the reference genome, where the second mate pair aligns. As with RNEXT, this field is only applicable when processing paired-end sequencing data. A value of 0 indicates that information is not available.   | the <i>E. coli</i> simulated data is single read sequencing data, and all RNEXT values are therefore set to 0 (no information available). |
| TLEN       | Template Length. Only applicable for paired-end sequencing data, TLEN is the size of the original DNA or RNA fragment, determined by examining both of the <b>paired-mates</b> , and counting bases from the left-most aligned base to the right-most aligned base. A value of 0 indicates that TLEN information is not available.  | the <i>E. coli</i> simulated data is single read sequencing data, and all RNEXT values are therefore set to 0 (no information available). |
| SEQ        | the raw sequence, as originally defined in the FASTQ file.  | CCAGGGCAGTGGCAGGT...  |
| QUAL       | The <b>Phred quality score</b> for each base, as originally defined in the FASTQ file.  | .   |

# Converting SAM to BAM

- Download Samtools (<http://www.htslib.org/>)
- To convert from SAM to BAM, use the SAMtools view command:

```
$ samtools view -b -S -o sim_reads_1M_aligned.bam  
sim_reads_1M_aligned.sam
```

- -b: indicates that the output is BAM.
  - -S: indicates that the input is SAM.
  - -o: specifies the name of the output file.
- BAM files are stored in a compressed, binary format, and cannot be viewed directly. However, you can use the same view command to display all alignments. For example, running:

```
$ samtools view sim_reads_1Ms_aligned.bam | more
```

# Sorting and Indexing

- The next step is to sort and index the BAM file.

```
$ samtools sort sim_reads_1M_aligned.bam  
sim_reads_1M_aligned.sorted
```

- Once you have sorted your BAM file, you can then index it.

```
$ samtools index sim_reads_1M_aligned.sorted.bam
```

# Identifying Genomic Variants

- The first step is to use the SAMtools mpileup command to calculate the **genotype likelihoods** supported by the aligned reads in our sample:

```
$ samtools mpileup -g -f NC_008253.fna  
sim_reads_1M_aligned.sorted.bam >  
sim_reads_1M_variants.bcf
```

-g: directs SAMtools to output genotype likelihoods in the **binary call format (BCF)**. This is a compressed binary format.

- f: directs SAMtools to use the specified reference genome. A reference genome must be specified, and here we specify the reference genome for *E. coli*.

# Identifying Genomic Variants

- The second step is to use bcftools:

```
$ bcftools call -c -v sim_reads_1M_variants.bcf >  
sim_reads_1M_variants.vcf
```

- c: directs bcftools to call SNPs.
- v: directs bcftools to only output potential variants

# Visualizing Reads

- For the final step, we will use the SAMtools tview command to view our simulated reads and visually compare them to the *E. coli* reference genome. To do so, run:

```
$ samtools tview sim_reads_aligned.sorted.bam  
NC_008253.fna
```

# GATK

- <https://www.broadinstitute.org/gatk/>

The Genome Analysis Toolkit or GATK is a software package for analysis of high-throughput sequencing data, developed by the [Data Science and Data Engineering](#) group at the [Broad Institute](#). The toolkit offers a wide variety of tools, with a primary focus on variant discovery and genotyping as well as strong emphasis on data quality assurance. Its robust architecture, powerful processing engine and high-performance computing features make it capable of taking on projects of any size.

[Learn more »](#)



- <https://www.broadinstitute.org/gatk/guide/best-practices>
- <http://broadinstitute.github.io/picard/>
- <https://broadinstitute.github.io/picard/command-line-overview.html>

# Picard

- <http://broadinstitute.github.io/picard/>
- <https://broadinstitute.github.io/picard/command-line-overview.html>

The screenshot shows the GitHub project page for Picard. At the top left is the Picard logo with a green "build passing" badge. To the right are four buttons: "Latest Release", "Download ZIP File", "Download TAR Ball", and "View On GitHub". Below these buttons is a brief description: "A set of tools (in Java) for working with next generation sequencing data in the BAM (<http://samtools.sourceforge.net>) format." A large text box below contains more detailed information: "A set of Java command line tools for manipulating high-throughput sequencing data (HTS) data and formats. Picard is implemented using the HTSJDK Java library HTSJDK, supporting accessing of common file formats, such as SAM and VCF, used for high-throughput sequencing data."

- [https://software.broadinstitute.org/gatk/events/slides/1506/GATKwr8-A-3-GATK Best Practices and Broad pipelines.pdf](https://software.broadinstitute.org/gatk/events/slides/1506/GATKwr8-A-3-GATK_Best_Practices_and_Broad_pipelines.pdf)

# IGV

- <https://www.broadinstitute.org/igv/>

The screenshot shows the homepage of the Integrative Genomics Viewer (IGV) website. At the top left is the IGV logo and a search bar. The main header reads "Home" and "Integrative Genomics Viewer". Below the header is a large image showing a screenshot of the IGV software interface displaying genomic tracks. To the left of the main content area is a sidebar with links to "Home", "Downloads", "Documents", "Hosted Genomes", "FAQ", "IGV User Guide", "File Formats", "Release Notes", "IGV for iPad", "Credits", and "Contact". There is also a search bar and a "BROAD INSTITUTE" logo with the copyright notice "© 2013 Broad Institute". The main content area is divided into sections: "Overview", "Downloads", "Funding", and "Citing IGV". The "Overview" section contains a brief description of IGV and its capabilities. The "Downloads" section includes a download button and instructions for registration. The "Funding" section lists the funding sources for IGV. The "Citing IGV" section provides citation information for the tool.

**Home**

# Integrative Genomics Viewer

**Overview**

The Integrative Genomics Viewer (IGV) is a high-performance visualization tool for interactive exploration of large, integrated genomic datasets. It supports a wide variety of data types, including array-based and next-generation sequence data, and genomic annotations.

**Citing IGV**

To cite your use of IGV in your publication:

James T. Robinson, Helga Thorvaldsdóttir, Wendy Winckler, Mitchell Guttman, Eric S. Lander, Gad Getz, Jill P. Mesirov. [Integrative Genomics Viewer](#). *Nature Biotechnology* 29, 24–26 (2011).

Helga Thorvaldsdóttir, James T. Robinson, Jill P. Mesirov. [Integrative Genomics Viewer \(IGV\): high-performance genomics data visualization and exploration](#). *Briefings in Bioinformatics* 14, 178–192 (2013).

**Downloads**

Please [register](#) to download IGV. After registering, you can log in at any time using your email address.

**Funding**

Development of IGV is made possible by funding from the [National Cancer Institute](#), the [National Institute of General Medical Sciences](#) of the [National Institutes of Health](#), and the [Starr Cancer Consortium](#).

IGV participates in the [GenomeSpace](#) initiative, which is funded by the [National Human Genome Research Institute](#).