



BIMM 143

Genome Informatics I

Lecture 14

Barry Grant
UC San Diego

<http://thegrantlab.org/bimm143>

TODAYS MENU:

▶ **What is a Genome?**

- Genome sequencing and the Human genome project

▶ **What can we do with a Genome?**

- Compare, model, mine and edit

▶ **Modern Genome Sequencing**

- 1st, 2nd and 3rd generation sequencing

▶ **Workflow for NGS**

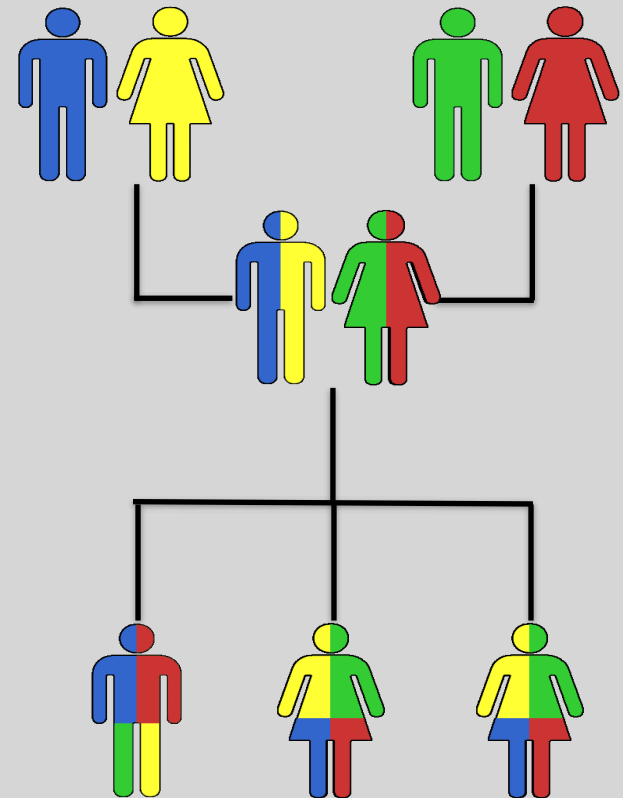
- RNA-Sequencing and Discovering variation

Genetics and Genomics

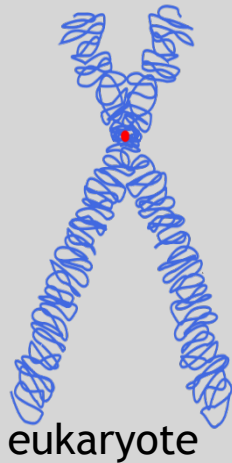
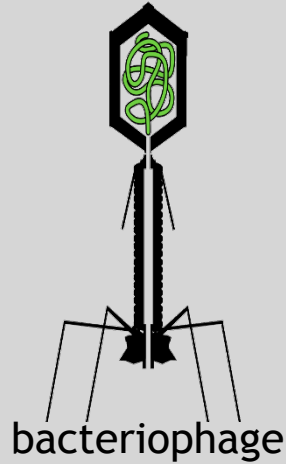
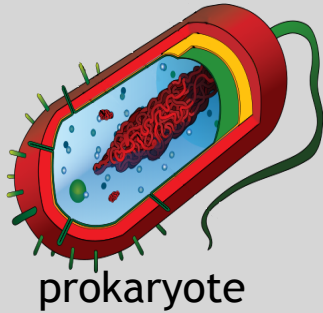
- **Genetics** is primarily the study of individual genes, mutations within those genes, and their inheritance patterns in order to understand specific traits.
- **Genomics** expands upon classical genetics and considers aspects of the entire genome, typically using computer aided approaches.

What is a Genome?

The total genetic material of an organism by which individual traits are encoded, controlled, and ultimately passed on to future generations

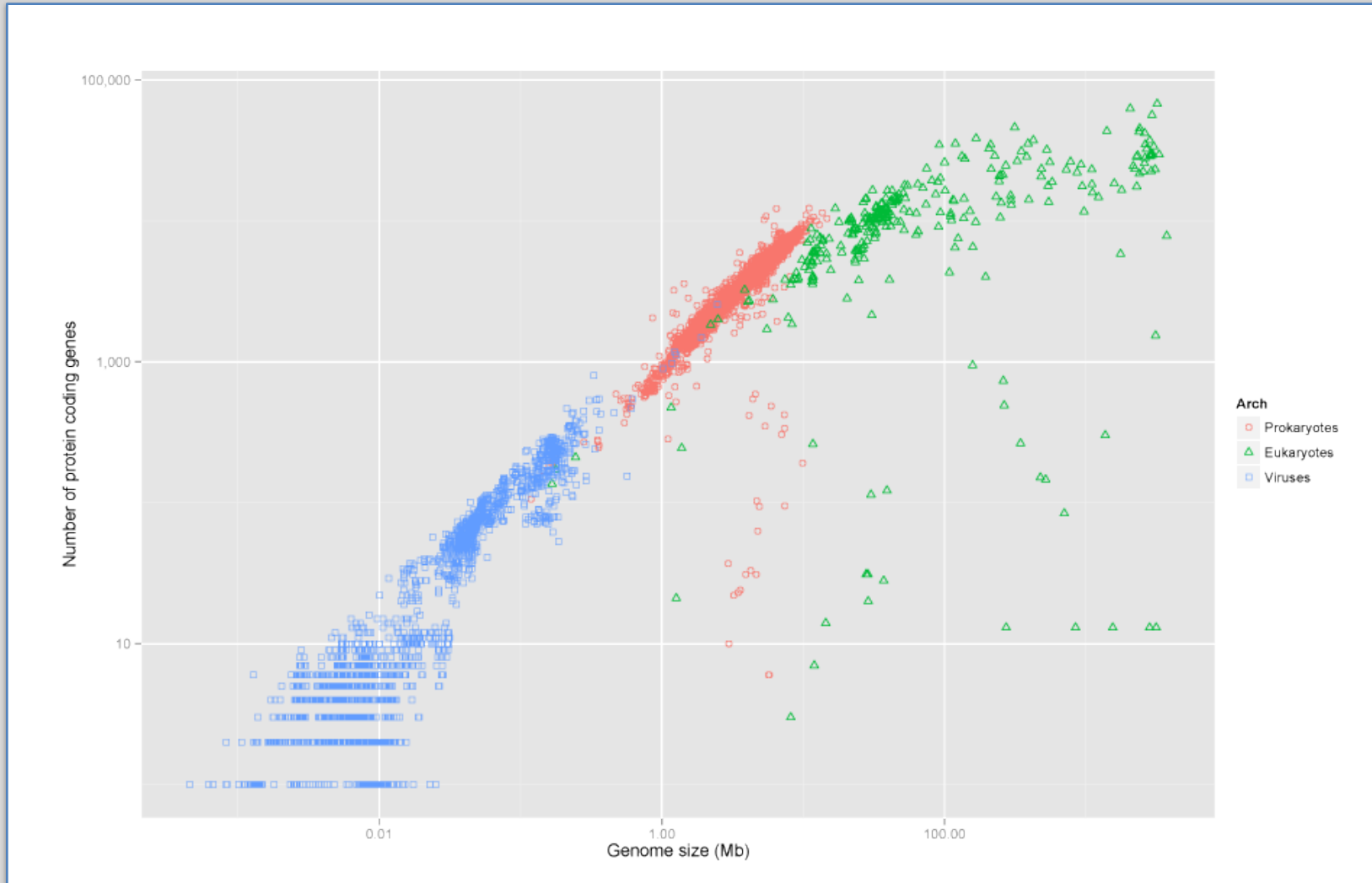


Genomes come in many shapes



- Primarily DNA, but can be RNA in the case of some viruses
- Some genomes are circular, others linear
- Can be organized into discrete units (chromosomes) or freestanding molecules (plasmids)

Genomes come in many sizes




Genome Databases

NCBI Genome:

<http://www.ncbi.nlm.nih.gov/genome>

NCBI Resources How To Sign in to NCBI

Genome Search Limits Advanced Help



Genome

This resource organizes information on genomes including sequences, maps, chromosomes, assemblies, and annotations.

Using Genome

- [Help](#)
- [Browse by Organism](#)
- [Download / FTP](#)
- [Download FAQ](#)
- [Submit a genome](#)

Genome Tools

- [BLAST the Human Genome](#)
- [Microbial Nucleotide BLAST](#)
- [TaxPlot \(3-way Genome Comparison\)](#)

Custom resources

- [Human Genome](#)
- [Microbes](#)
- [Organelles](#)
- [Viruses](#)
- [Prokaryotic reference genomes](#)

Genome Annotation and Analysis

- [Eukaryotic Genome Annotation](#)
- [Prokaryotic Genome Annotation](#)
- [PASC \(Pairwise Sequence Comparison\)](#)

Other Resources

- [Assembly](#)
- [BioProject](#)
- [BioSample](#)
- [Map Viewer](#)
- [Protein Clusters](#)


External Resources

- [GOLD - Genomes Online Database](#)
- [Ensembl Genome Browser](#)
- [Bacteria Genomes at Sanger](#)
- [Large-Scale Genome Sequencing \(NHGRI\)](#)

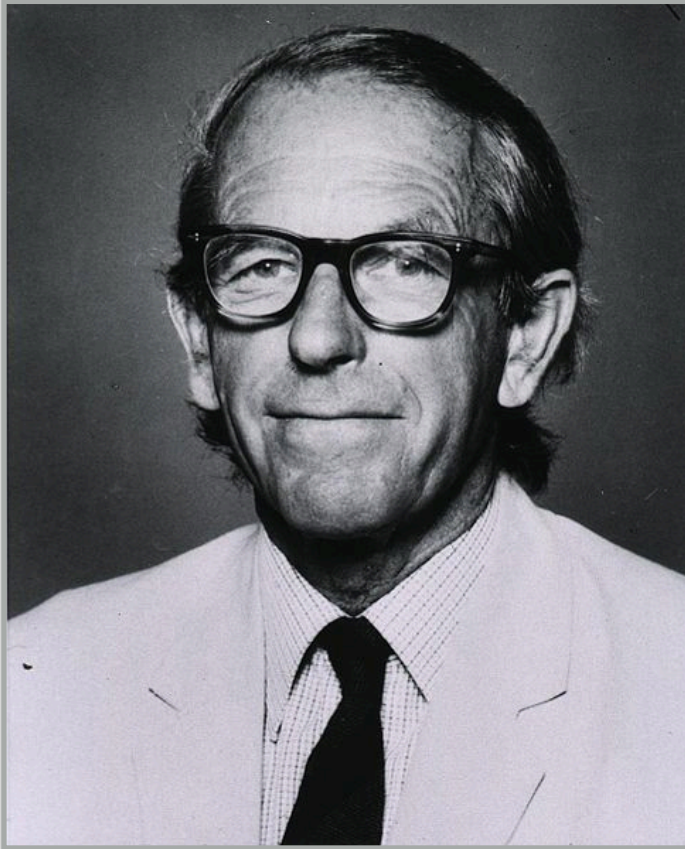
You are here: NCBI > Genomes & Maps > Genome Write to the Help Desk

GETTING STARTED <ul style="list-style-type: none">NCBI EducationNCBI Help ManualNCBI HandbookTraining & Tutorials	RESOURCES <ul style="list-style-type: none">Chemicals & BioassaysData & SoftwareDNA & RNADomains & StructuresGenes & ExpressionGenetics & MedicineGenomes & MapsHomologyLiteratureProteinsSequence AnalysisTaxonomyTraining & TutorialsVariation	POPULAR <ul style="list-style-type: none">PubMedBookshelfPubMed CentralPubMed HealthBLASTNucleotideGenomeSNPGeneProteinPubChem	FEATURED <ul style="list-style-type: none">Genetic Testing RegistryPubMed HealthGenBankReference SequencesGene Expression OmnibusMap ViewerHuman GenomeMouse GenomeInfluenza VirusPrimer-BLASTSequence Read Archive	NCBI INFORMATION <ul style="list-style-type: none">About NCBIResearch at NCBINCBI NewsNCBI FTP SiteNCBI on FacebookNCBI on TwitterNCBI on YouTube
---	--	---	--	--

Copyright | Disclaimer | Privacy | Browsers | Accessibility | Contact
National Center for Biotechnology Information, U.S. National Library of Medicine
8600 Rockville Pike, Bethesda MD, 20894 USA

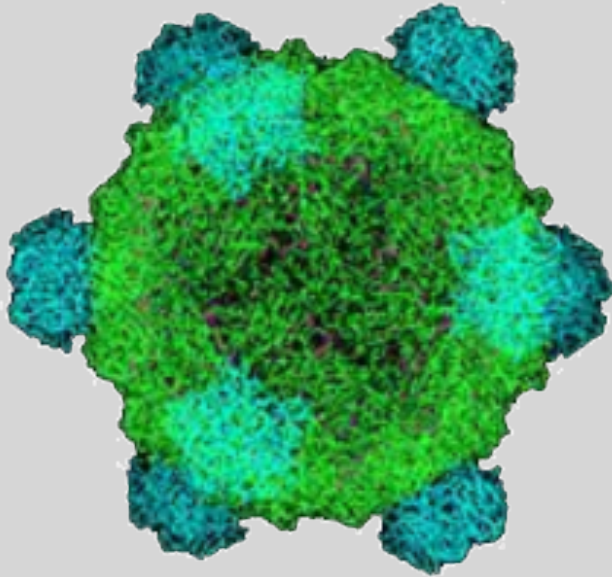


Early Genome Sequencing



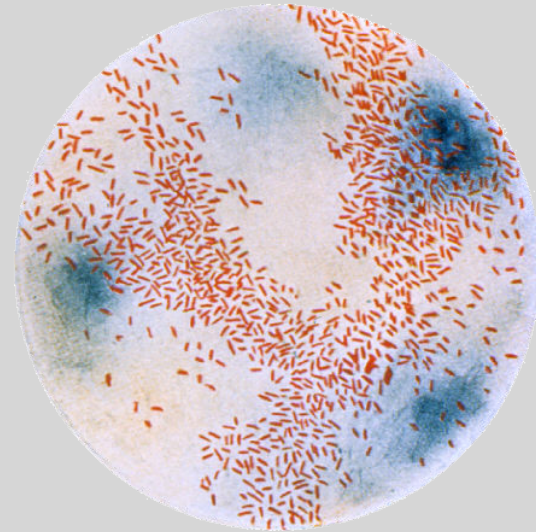
- Chain-termination “Sanger” sequencing was developed in 1977 by Frederick Sanger, colloquially referred to as the “Father of Genomics”
- Sequence reads were typically 750-1000 base pairs in length with an error rate of $\sim 1 / 10000$ bases

The First Sequenced Genomes



Bacteriophage ϕ -X174

- Completed in 1977
- 5,386 base pairs, ssDNA
- 11 genes



Haemophilus influenzae

- Completed in 1995
- 1,830,140 base pairs, dsDNA
- 1740 genes

The Human Genome Project

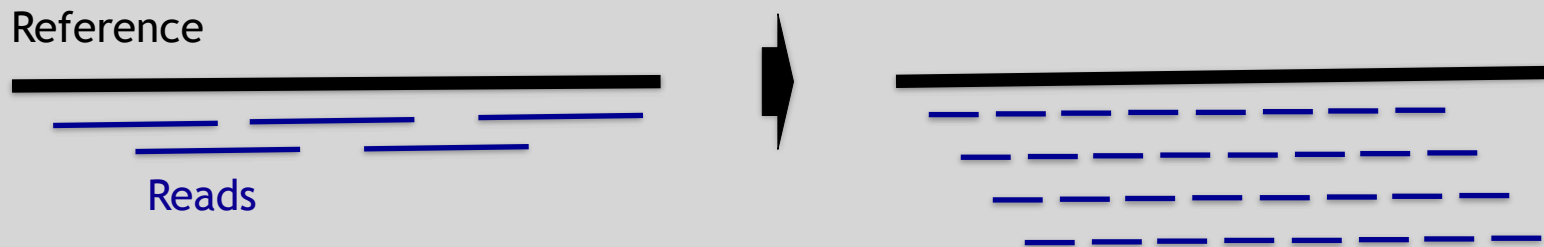
- The Human Genome Project (HGP) was an international, public consortium that began in 1990
 - Initiated by James Watson
 - Primarily led by Francis Collins
 - Eventual Cost: \$2.7 Billion
- Celera Genomics was a private corporation that started in 1998
 - Headed by Craig Venter
 - Eventual Cost: \$300 Million
- Both initiatives released initial drafts of the human genome in 2001
 - ~3.2 Billion base pairs, dsDNA
 - 22 autosomes, 2 sex chromosomes
 - ~20,000 genes



HHMI

Modern Genome Sequencing

- Next Generation Sequencing (NGS) technologies have resulted in a paradigm shift from long reads at low coverage to short reads at high coverage
- This provides numerous opportunities for new and expanded genomic applications



Rapid progress of genome sequencing

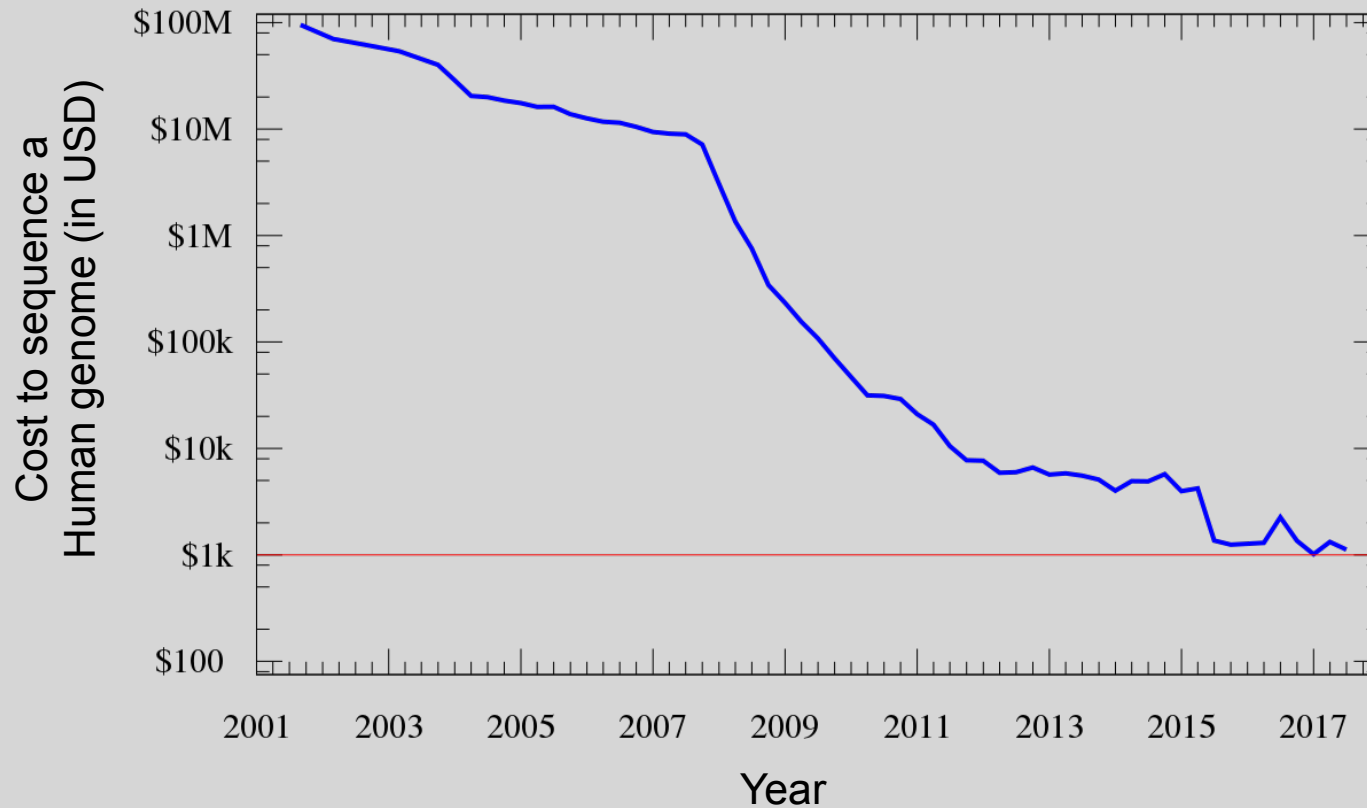


Image source: https://en.wikipedia.org/wiki/Carlson_curve

Rapid progress of genome sequencing

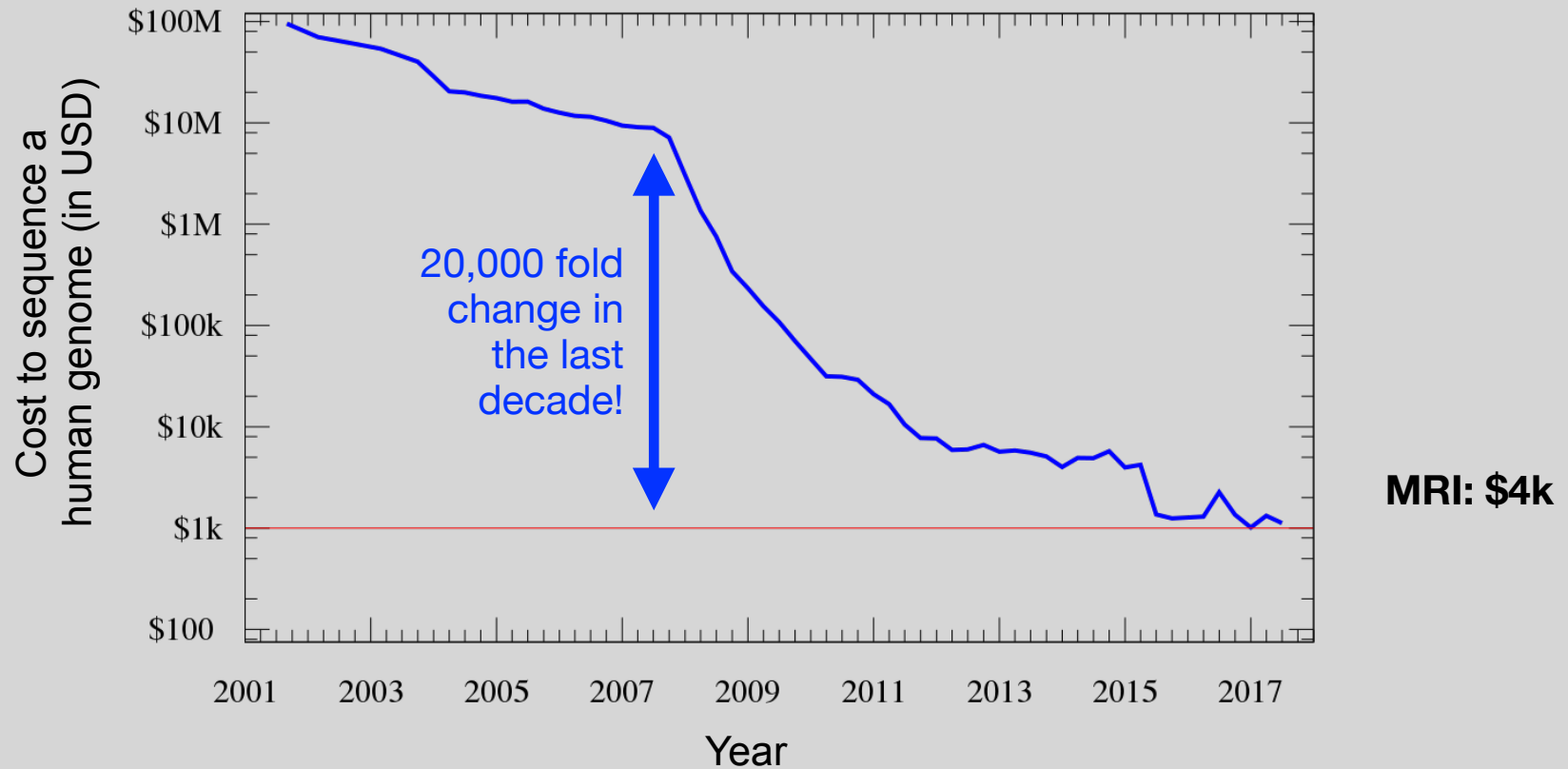
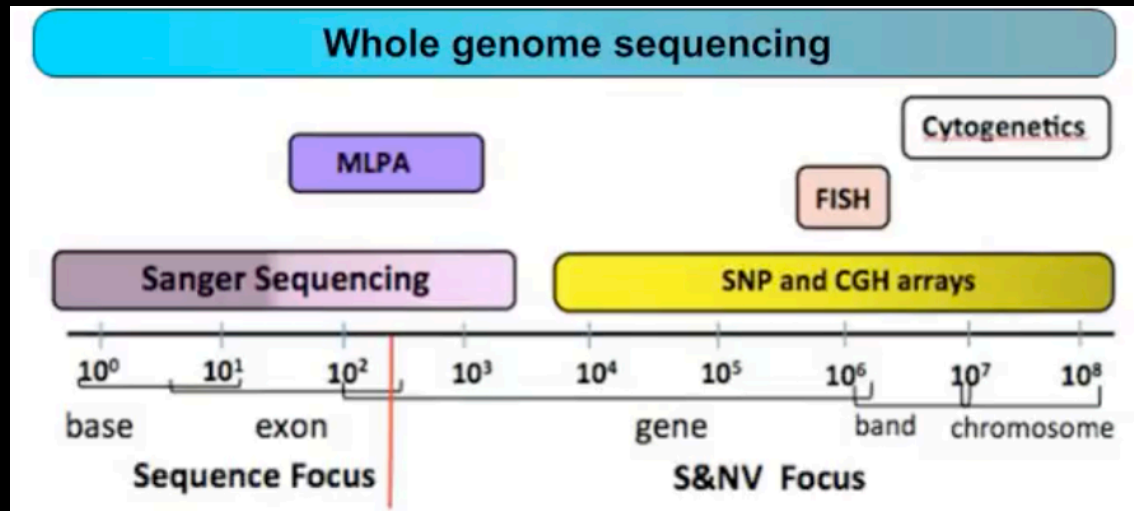


Image source: https://en.wikipedia.org/wiki/Carlson_curve

Whole genome sequencing transforms genetic testing



- 1000s of single gene tests
- Structural and copy number variation tests
- Permits hypothesis free diagnosis

Major impact areas for genomic medicine

- **Cancer**: Identification of driver mutations and drugable variants, Molecular stratification to guide and monitor treatment, Identification of tumor specific variants for personalized immunotherapy approaches (precision medicine).
- **Genetic disease diagnose**: Rare, inherited and so-called ‘mystery’ disease diagnose.
- **Health management**: Predisposition testing for complex diseases (e.g. cardiac disease, diabetes and others), optimization and avoidance of adverse drug reactions.
- **Health data analytics**: Incorporating genomic data with additional health data for improved healthcare delivery.

Goals of Cancer Genome Research

- Identify changes in the genomes of tumors that drive cancer progression
- Identify new targets for therapy
- Select drugs based on the genomics of the tumor
- Provide early cancer detection and treatment response monitoring
- Utilize cancer specific mutations to derive neoantigen immunotherapy approaches



What can go wrong in cancer genomes?

Type of change	Some common technology to study changes
DNA mutations	WGS, WXS
DNA structural variations	WGS
Copy number variation (CNV)	CGH array, SNP array, WGS
DNA methylation	Methylation array, RRBS, WGBS
mRNA expression changes	mRNA expression array, RNA-seq
miRNA expression changes	miRNA expression array, miRNA-seq
<i>Protein expression</i>	Protein arrays, mass spectrometry

WGS = whole genome sequencing, WXS = whole exome sequencing

RRBS = reduced representation bisulfite sequencing, WGBS = whole genome bisulfite sequencing

DNA Sequencing Concepts

- **Sequencing by Synthesis:** Uses a polymerase to incorporate and assess nucleotides to a primer sequence
 - 1 nucleotide at a time
- **Sequencing by Ligation:** Uses a ligase to attach hybridized sequences to a primer sequence
 - 1 or more nucleotides at a time (e.g. dibase)

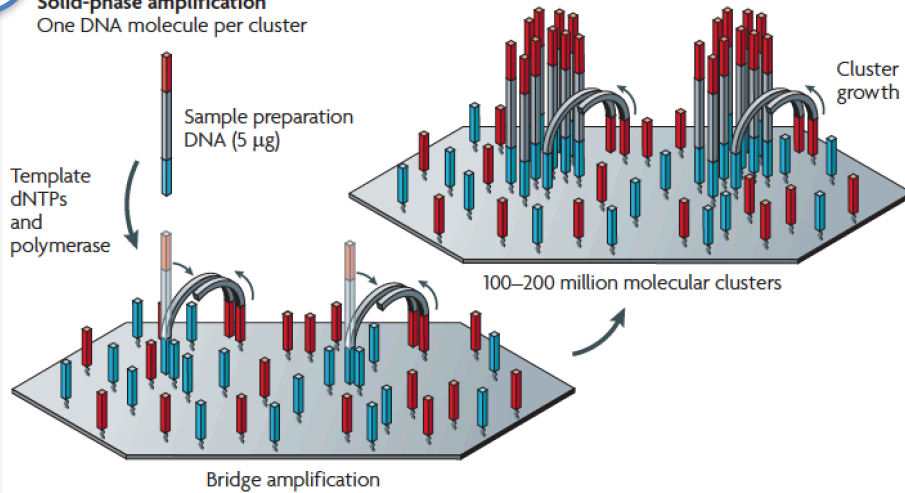
Modern NGS Sequencing Platforms

	Roche/454	Life Technologies SOLiD	Illumina Hi-Seq 2000
Library amplification method	emPCR* on bead surface	emPCR* on bead surface	Enzymatic amplification on glass surface
Sequencing method	Polymerase-mediated incorporation of unlabelled nucleotides	Ligase-mediated addition of 2-base encoded fluorescent oligonucleotides	Polymerase-mediated incorporation of end-blocked fluorescent nucleotides
Detection method	Light emitted from secondary reactions initiated by release of PPI	Fluorescent emission from ligated dye-labelled oligonucleotides	Fluorescent emission from incorporated dye-labelled nucleotides
Post incorporation method	NA (unlabelled nucleotides are added in base-specific fashion, followed by detection)	Chemical cleavage removes fluorescent dye and 3' end of oligonucleotide	Chemical cleavage of fluorescent dye and 3' blocking group
Error model	Substitution errors rare, insertion/deletion errors at homopolymers	End of read substitution errors	End of read substitution errors
Read length (fragment/paired end)	400 bp/variable length mate pairs	75 bp/50+25 bp	150 bp/100+100 bp

Illumina - Reversible terminators

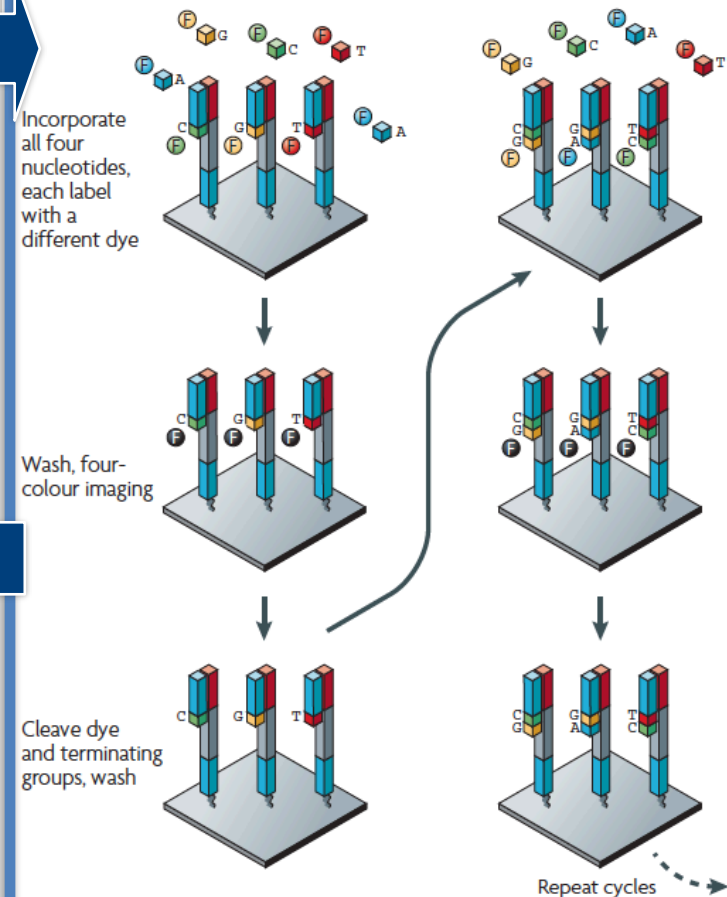
1 Enzymatic amplification on glass surface

Illumina/Solexa
Solid-phase amplification
One DNA molecule per cluster

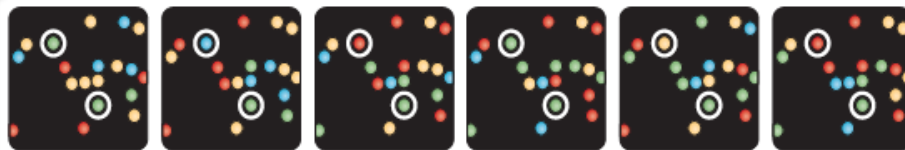


2 Polymerase-mediated incorporation of end blocked fluorescent nucleotides

Illumina/Solexa — Reversible terminators



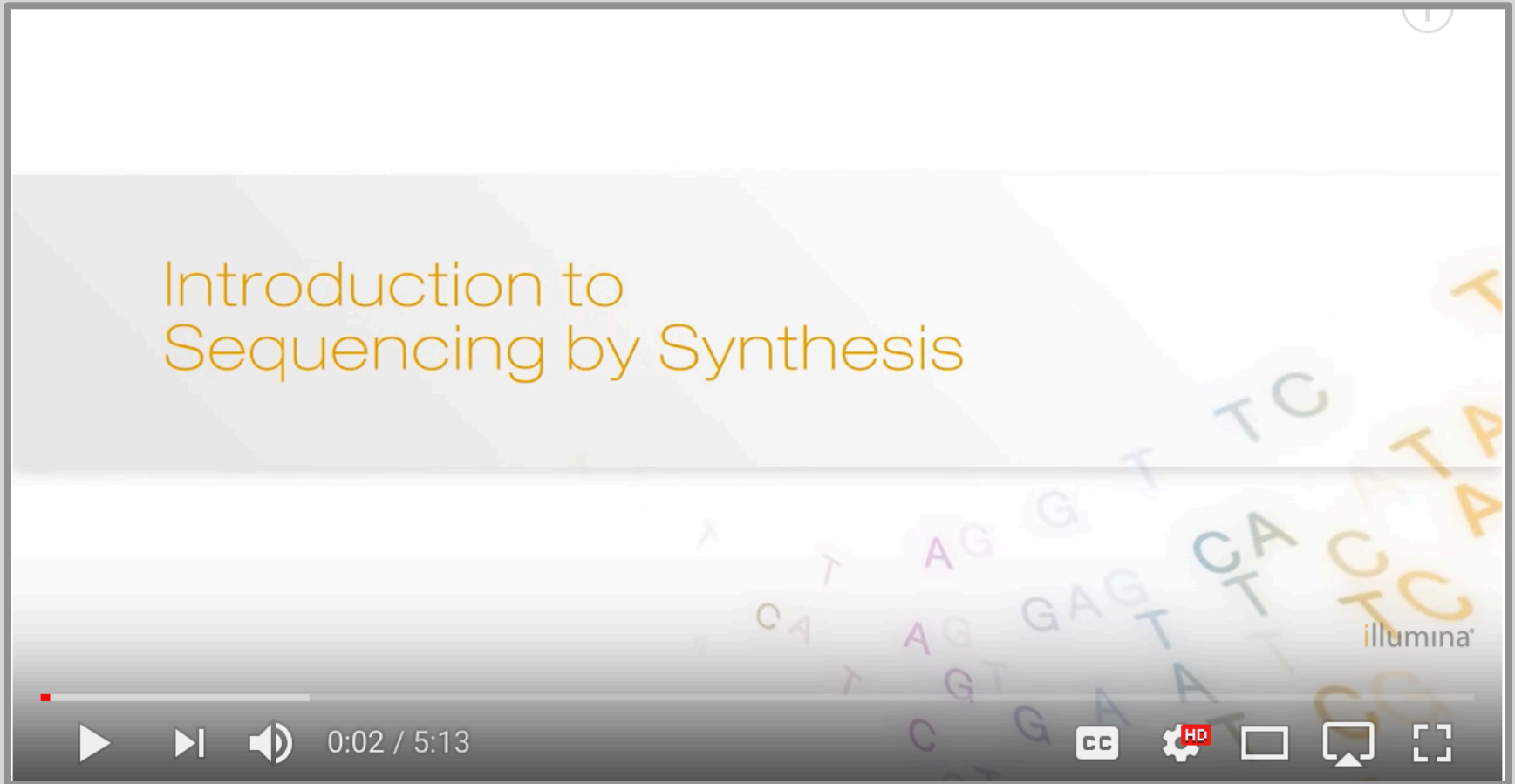
3 Fluorescent emission from incorporated dye-labeled nucleotides



Top: CATCGT
Bottom: CCCCCC

Cleave dye & blocking group, repeat...

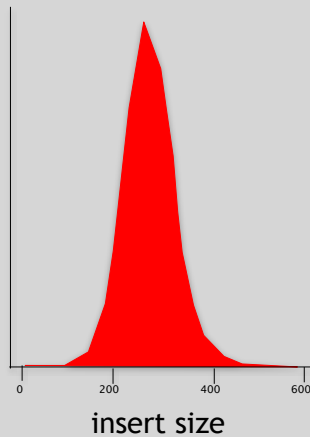
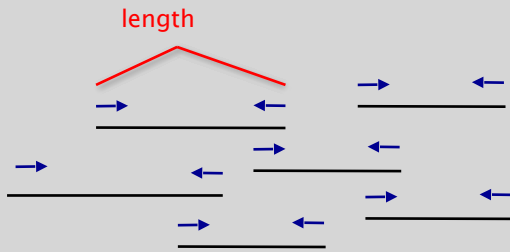
Illumina Sequencing - Video



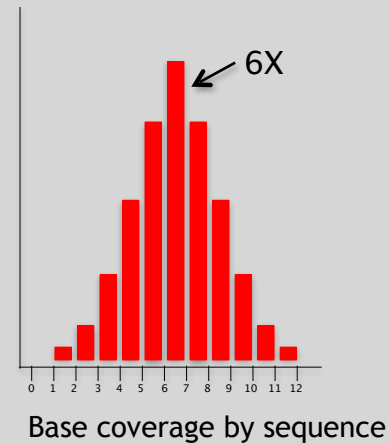
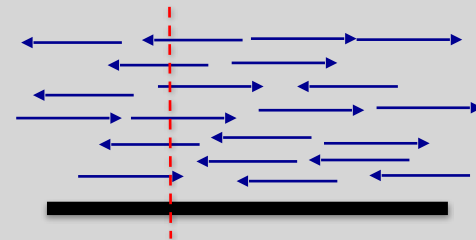
https://www.youtube.com/watch?src_vid=womKfikWlxM&v=fCd6B5HRaZ8

NGS Sequencing Terminology

Insert Size



Sequence Coverage



Summary: “Generations” of DNA Sequencing

	First generation	Second generation ^a	Third generation ^a
Fundamental technology	Size-separation of specifically end-labeled DNA fragments, produced by SBS or degradation	Wash-and-scan SBS	SBS, by degradation, or direct physical inspection of the DNA molecule
Resolution	Averaged across many copies of the DNA molecule being sequenced	Averaged across many copies of the DNA molecule being sequenced	Single-molecule resolution
Current raw read accuracy	High	High	Moderate
Current read length	Moderate (800–1000 bp)	Short, generally much shorter than Sanger sequencing	Long, 1000 bp and longer in commercial systems
Current throughput	Low	High	Moderate
Current cost	High cost per base Low cost per run	Low cost per base High cost per run	Low-to-moderate cost per base Low cost per run
RNA-sequencing method	cDNA sequencing	cDNA sequencing	Direct RNA sequencing and cDNA sequencing
Time from start of sequencing reaction to result	Hours	Days	Hours
Sample preparation	Moderately complex, PCR amplification not required	Complex, PCR amplification required	Ranges from complex to very simple depending on technology
Data analysis	Routine	Complex because of large data volumes and because short reads complicate assembly and alignment algorithms	Complex because of large data volumes and because technologies yield new types of information and new signal processing challenges
Primary results	Base calls with quality values	Base calls with quality values	Base calls with quality values, potentially other base information such as kinetics

Third Generation Sequencing

- Currently in active development
- Hard to define what “3rd” generation means
- Typical characteristics:
 - Long (1,000bp+) sequence reads
 - Single molecule (no amplification step)
 - Often associated with nanopore technology
 - But not necessarily!

The first direct RNA sequencing by nanopore

Side-Note:

- For example this new nanopore sequencing method was just published!

<https://www.nature.com/articles/nmeth.4577>

- "Sequencing the RNA in a biological sample can unlock a wealth of information, including the identity of bacteria and viruses, the nuances of alternative splicing or the transcriptional state of organisms. However, current methods have limitations due to short read lengths and reverse transcription or amplification biases. Here we demonstrate nanopore direct RNA-seq, a highly parallel, real-time, single-molecule method that circumvents reverse transcription or amplification steps."

SeqAnswers Wiki

Side-Note:

A good repository of analysis software can be found at <http://seqanswers.com/wiki/Software/list>

Page [Discussion](#) [Read](#) [View source](#) [View history](#) [Log in](#)

Software/list

[< Software](#)

Below is (one of many possible) dynamic tables of software data, created from pages in the wiki. To add a package to the list, use the following form:

CSV
JSON

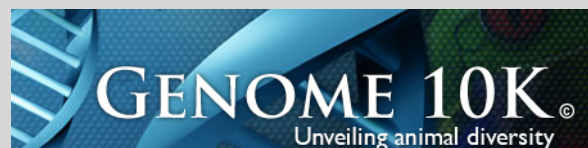
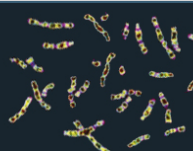
Name	Summary	Bio Tags	Meth Tags	Features	Language	Licence	OS
4peaks	Allows viewing sequencing trace files, motif searching trimming, BLAST and exporting sequences.	Sequencing	Sequence analysis			Freeware	Mac OS X
AB Large Indel Tool	Identifies deviations in clone insert size that indicate intra-chromosomal structural variations compared to a reference genome.	InDel discovery Sequencing	Mapping		Perl	GPL	Linux 64
AB Small Indel Tool	The SOLID™ Small Indel Tool processes the indel evidences found in the pairing step of the SOLID™ System Analysis Pipeline Tool (Corona Lite).	InDel discovery Sequencing	Mapping Alignment		Perl C++	GPL	Linux 64
ABBA	Assembly Boosted By Amino acid sequence is a comparative gene assembler, which uses amino acid sequences from predicted proteins to help build a better assembly	Genomic Assembly	Assembly Scaffolding			Artistic License	Linux
ABMapper	Maps RNA-Seq reads to target genome considering possible multiple mapping locations and splice junctions	Genomics Transcriptomics	Mapping Alignment		C++ Perl	GPLv3	Linux
ABYSS	ABYSS is a de novo sequence assembler designed for short reads and large genomes.	De-novo assembly	Assembly De Bruijn graph	MPI OpenMP	C++	Free for academic use	POSIX Linux Mac OS X
Adaptor Removal	Removes adaptor fragments from raw short read	General	Adaptor Removal	Trimming	Java	Custom Licence	Linux 64

**What can we do with all
this sequence information?**

Population Scale Analysis

We can now begin to assess genetic differences on a very large scale, both as naturally occurring variation in human and non-human populations as well somatically within tumors

1000 Genomes
A Deep Catalog of Human Genetic Variation



The Cancer Genome Atlas



*Understanding genomics
to improve cancer care*

The 100,000 Genomes Project

Genomics England & Partners



<https://www.genomicsengland.co.uk/the-100000-genomes-project/>

“Variety’s the very spice of life”

-William Cowper, 1785

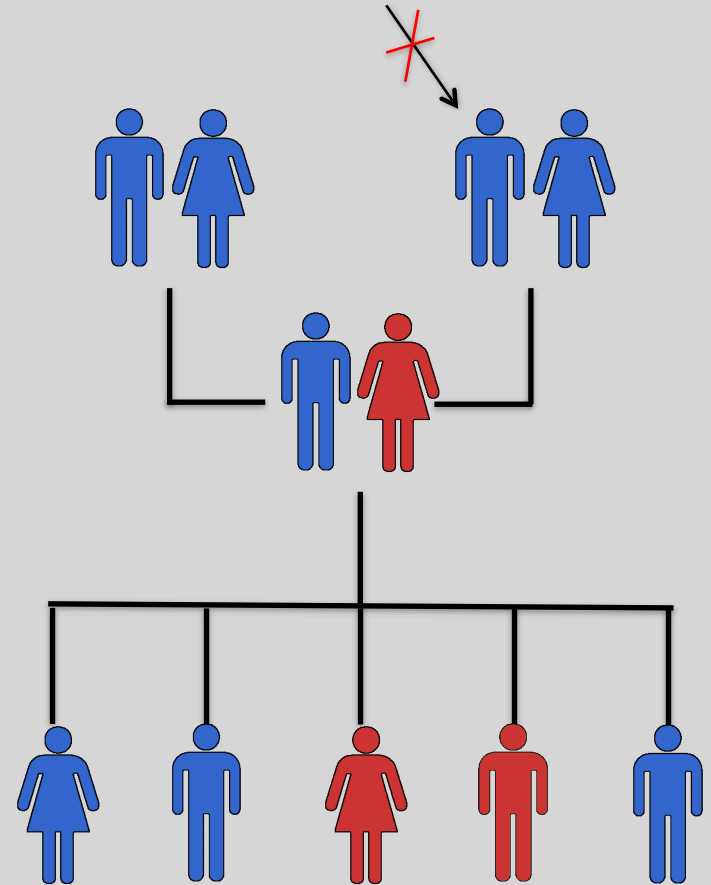
“Variation is the spice of life”

-Kruglyak & Nickerson, 2001

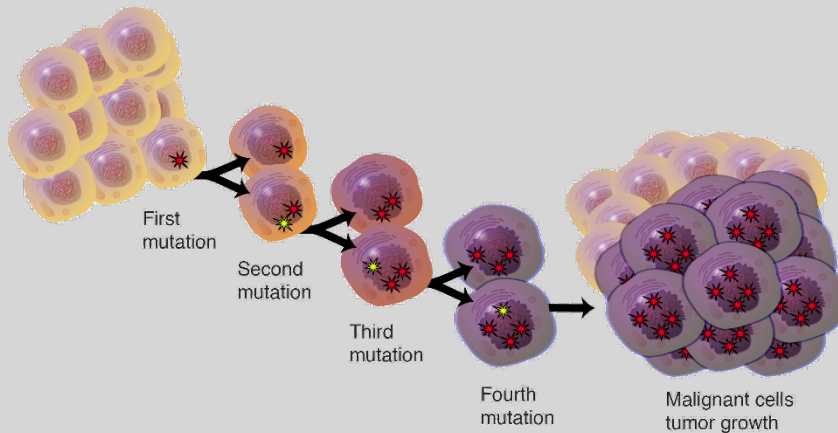
- While the sequencing of the human genome was a great milestone, the DNA from a single person is not representative of the millions of potential differences that can occur between individuals
- These unknown genetic variants could be the cause of many phenotypes such as differing morphology, susceptibility to disease, or be completely benign.

Germline Variation

- Mutations in the germline are passed along to offspring and are present in the DNA over every cell
- In animals, these typically occur in meiosis during gamete differentiation



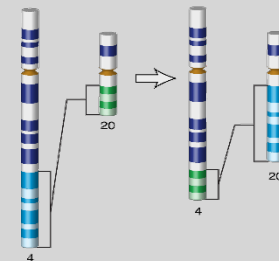
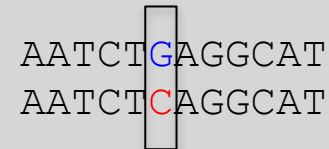
Somatic Variation



- Mutations in non-germline cells that are not passed along to offspring
- Can occur during mitosis or from the environment itself
- Are an integral part in tumor progression and evolution

Types of Genomic Variation

- **Single Nucleotide Polymorphisms (SNPs)** - mutations of one nucleotide to another
- **Insertion/Deletion Polymorphisms (INDELs)** - small mutations removing or adding one or more nucleotides at a particular locus
- **Structural Variation (SVs)** - medium to large sized rearrangements of chromosomal DNA



Differences Between Individuals

The average number of genetic differences in the germline between two random humans can be broken down as follows:

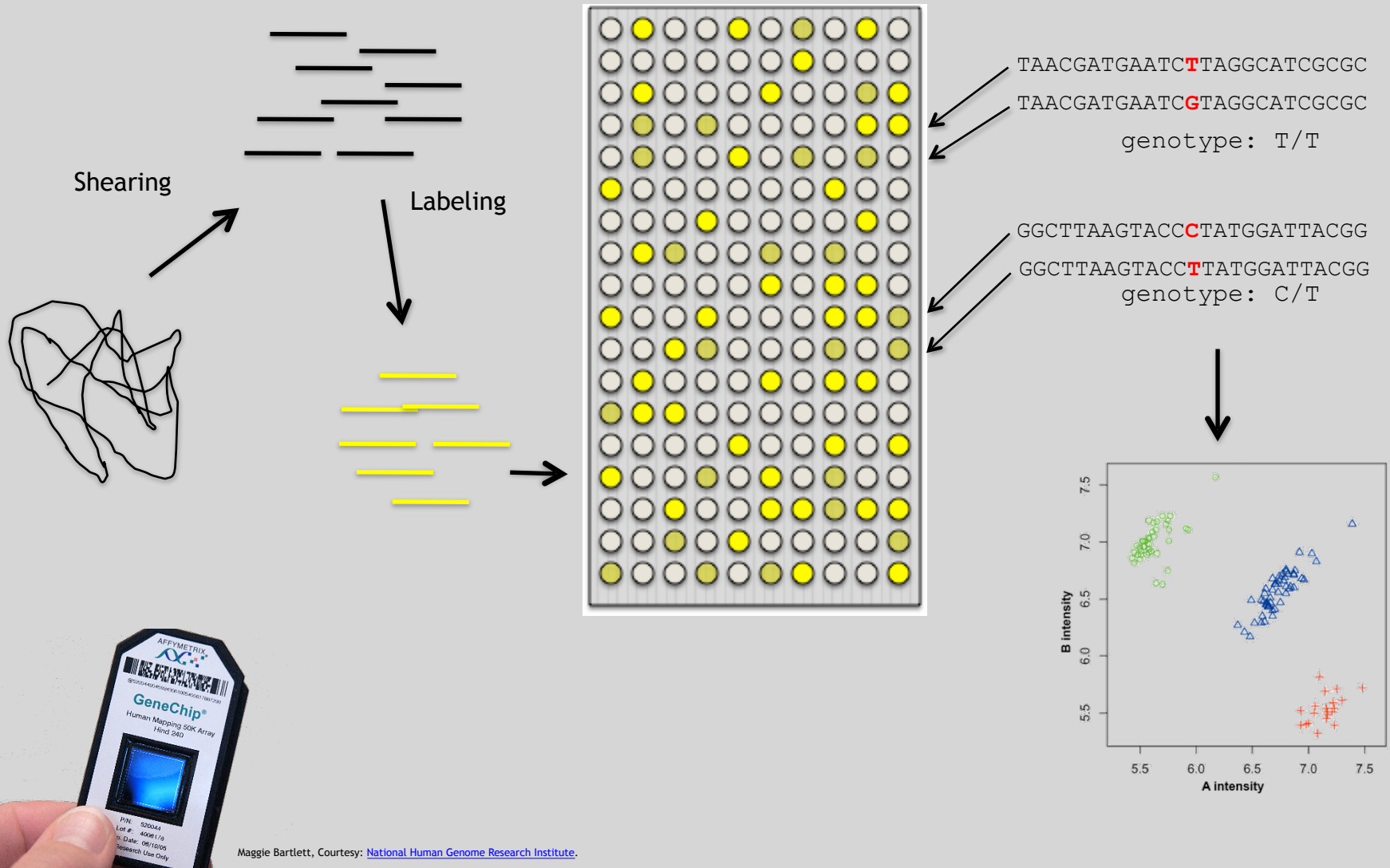
- 3,600,000 single nucleotide differences
- 344,000 small insertion and deletions
- 1,000 larger deletion and duplications

Numbers change depending on ancestry!

Genotyping Small Variants

- Once discovered, oligonucleotide probes can be generated with each individual allele of a variant of interest
- A large number can then be assessed simultaneously on microarrays to detect which combination of alleles is present in a sample

SNP Microarrays



Impact of Genetic Variation

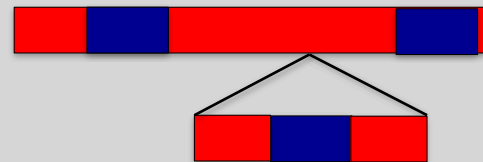
There are numerous ways genetic variation can exhibit functional effects

Premature stop codons



TAC -> TAA

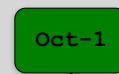
Gene or exon deletion



Frameshift mutation



TAC -> T-C



Transcription factor binding disruption



ATGCAAAT -> ATGCAGAT

Do it Yourself!

Hand-on time!

https://bioboot.github.io/bimm143_F18/lectures/#14

Sections **1** to **3** please (up to running Read Alignment)
See IP address on website for **your** Galaxy server

<http://uswest.ensembl.org/Help/View?id=140>

The image displays the Ensembl genome browser interface, illustrating three levels of genomic detail:

- Chromosome image:** Shows a chromosome with a highlighted **Region of interest** and **Haplotypes and patches**.
- Overview image:** Shows a detailed view of the **Region in detail**, including **Genes** and a **Gene or region of interest**. It includes a **Gene Legend** for merged Ensembl/Havana, processed transcript, pseudogene, and RNA gene.
- Zoomable Region image:** Shows a highly detailed view of the **Genome**, including **Transcripts (splice variants)** and **Genes**.

Callouts indicate **Change or add data tracks** for the configuration gear icons in the sidebar and the zoomable region image.

Access a jetstream galaxy instance!

Use assigned IP address

Do it Yourself!

Galaxy

149.165.169.186

Apps Gmail Seminars Atmosphere BGGN 213 · An intr...

Galaxy Analyze Data Workflow Shared Data Visualization Help User Using 12.3 MB

Tools

search tools

Get Data
Send Data
Collection Operations
Text Manipulation
Filter and Sort
Join, Subtract and Group
Convert Formats
Extract Features
Fetch Sequences
Fetch Alignments
Statistics
Graph/Display Data
FASTA manipulation
NGS: QC and manipulation
NGS: DeepTools
NGS: Mapping
Lastz map short reads against reference sequence
Map with Bowtie for Illumina
Map with BWA for Illumina
Map with BWA for SOLiD
Megablast compare short reads against htgs, nt, and wgs databases
Parse blast XML output
Map with BWA-MEM - map medium and long reads (> 100 bp) against reference genome
Map with BWA - map short reads (< 100 bp) against reference genome
Bowtie2 - map reads against reference genome
NGS: RNA Analysis

Bowtie2 - map reads against reference genome (Galaxy Version 2.2.6.2)

Options

Is this single or paired library
Single-end

FASTQ file
4: HG00109_2.fastq
Must be of datatype "fastqsanger"

Write unaligned reads (in fastq format) to separate file(s)
Yes No
--un/--un-conc; This triggers --un parameter for single reads and --un-conc for paired reads

Write aligned reads (in fastq format) to separate file(s)
Yes No
--al/--al-conc; This triggers --al parameter for single reads and --al-conc for paired reads

Will you select a reference genome from your history or use a built-in index?
Use a built-in genome index
Built-ins were indexed using default options. See `Indexes` section of help below

Select reference genome
Baboon (Papio anubis): papHam1
If your genome of interest is not listed, contact the Galaxy team

Set read groups information?
Do not set
Specifying read group information can greatly simplify your downstream analyses by allowing combining multiple datasets.

Select analysis mode
1: Default setting only

Do you want to use presets?
 No, just use defaults
 Very fast end-to-end (--very-fast)
 Fast end-to-end (--fast)
 Sensitive end-to-end (--sensitive)
 Very sensitive end-to-end (--very-sensitive)
 Very fast local (--very-fast-local)
 Fast local (--fast-local)
 Sensitive local (--sensitive-local)
 Very sensitive local (--very-sensitive-local)

Allow selecting among several preset parameter settings. Choosing between these will result in dramatic changes in runtime. See help below to understand effects of these presets.

History

search datasets

Unnamed history
22 shown, 2 deleted, 1 hidden
12.32 MB

25: htseq-count on data 18 and data 17 (no feature)
24: htseq-count on data 18 and data 17
23: Cufflinks on data 18 and data 16: Skipped Transcripts
21: Cufflinks on data 18 and data 16: assembled transcripts
20: Cufflinks on data 18 and data 16: transcript expression
19: Cufflinks on data 18 and data 16: gene expression

575 lines
format: tabular, database: hg19

```
cufflinks v2.2.1
cufflinks -q --no-update-check -l
300000 -F 0.100000 -j 0.150000 -p
6 -G /opt/galaxy/galaxy-
app/database/datasets/000/dataset_4
/opt/galaxy/galaxy-
app/database/datasets/000/dataset_4
```

1	2	3
tracking_id	class_code	nearest_ref_id
ZZEF1	-	-
CYB5D2	-	-
ANKFY1	-	-

Raw data usually in FASTQ format

```
@NS500177:196:HFTTTAFXX:1:11101:10916:1458 2:N:0:CGCGGCTG
ACACGACGATGAGGTGACAGTCACGGAGGATAAGATCAATGCCCTCATTAAAGCAGCCGGTGTAA
+
AAAAAEEEEEEEEEEEEEE//AEEEEEEEEEEEEEEEEEE/EE/<<EE/AEEFAEE///EEEEEEEEAEA<
```

1

2

3

4

Each sequencing “read” consists of 4 lines of data :

- 1 The first line (which always starts with ‘@’) is a unique ID for the sequence that follows
- 2 The second line contains the bases called for the sequenced fragment
- 3 The third line is always a “+” character
- 4 The fourth line contains the quality scores for each base in the sequenced fragment (these are ASCII encoded...)

ASCII Encoded Base Qualities

```
@NS500177:196:HFTTTAFXX:1:11101:10916:1458 2:N:0:CGCGGCTG
ACACGACGATGAGGTGACAGTCACGGAGGATAAGATCAATGCCCTCATTAAAGCAGCCGGTGTAA
+
AAAAAEEEEEEEEEEEEEE//AEEEEEEEEEEEEEEEE/EE/<<EE/AEEFAEE///EEEEEEEEAEA<
```

4

- Each sequence base has a corresponding numeric quality score encoded by a single ASCII character typically on the 4th line (see 4 above)
- ASCII characters represent integers between 0 and 127
- Printable ASCII characters range from 33 to 126
- Unfortunately there are 3 quality score formats that you may come across...

Interpreting Base Qualities in R

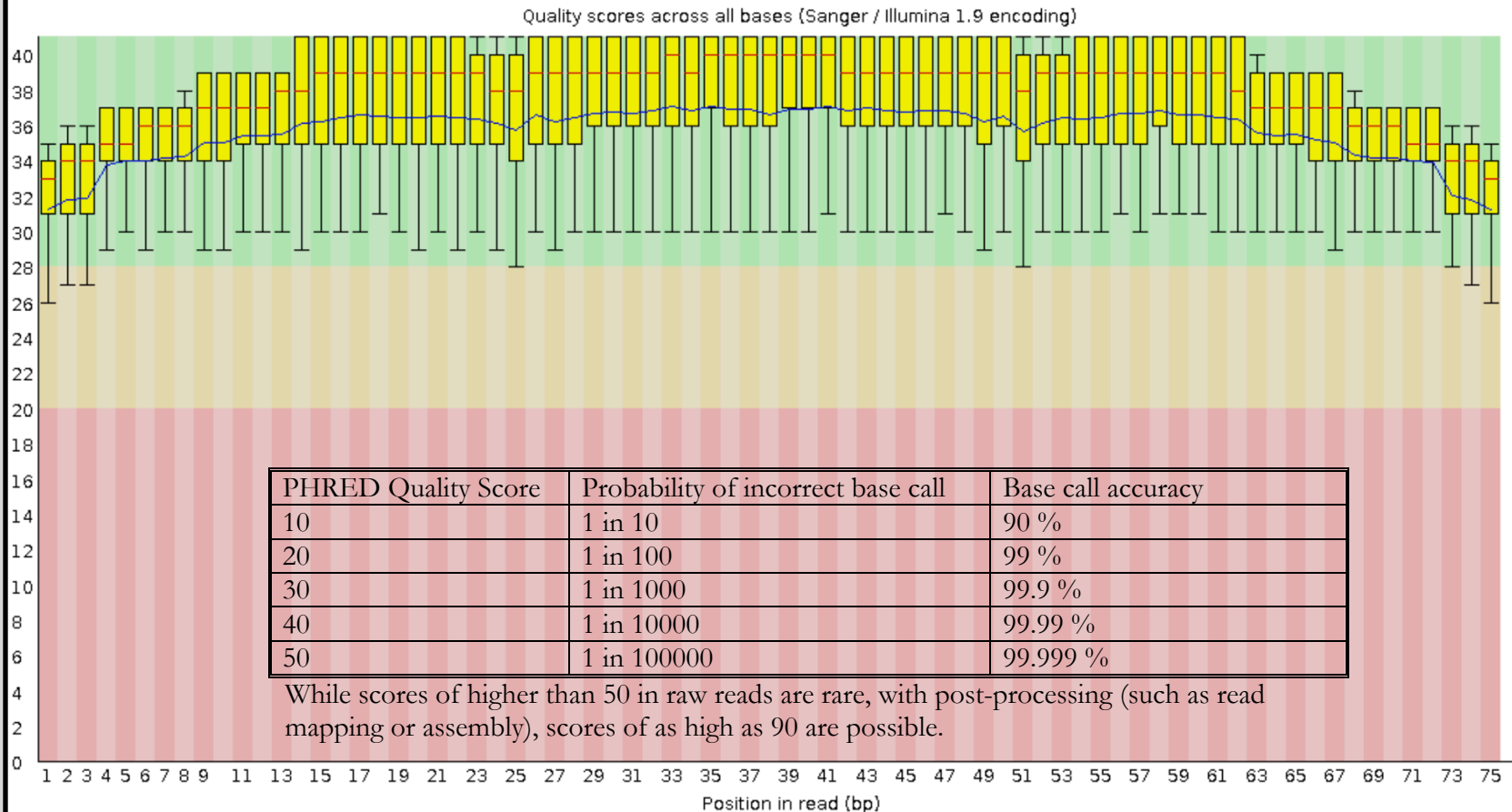
		ASCII Range	Offset	Score Range
Sanger, Illumina (Ver > 1.8)	fastqsanger	33-126	33	0-93
Solexa, Illumina (Ver < 1.3)	fastqsolexa	59-126	64	5-62
Illumina (Ver 1.3 -1.7)	fastqillumina	64-126	64	0-62

```
> library(seqinr)
> library(gtools)
> phred <- asc( s2c("DDDDCDEDCDDDDBBDDCC@") ) - 33
> phred
## D D D D C D E D C D D D D B B D D D C C @
## 35 35 35 35 34 35 36 35 34 35 35 35 35 33 33 35 35 35 34 34 31
> prob <- 10**(-phred/10)
```


FastQC Report



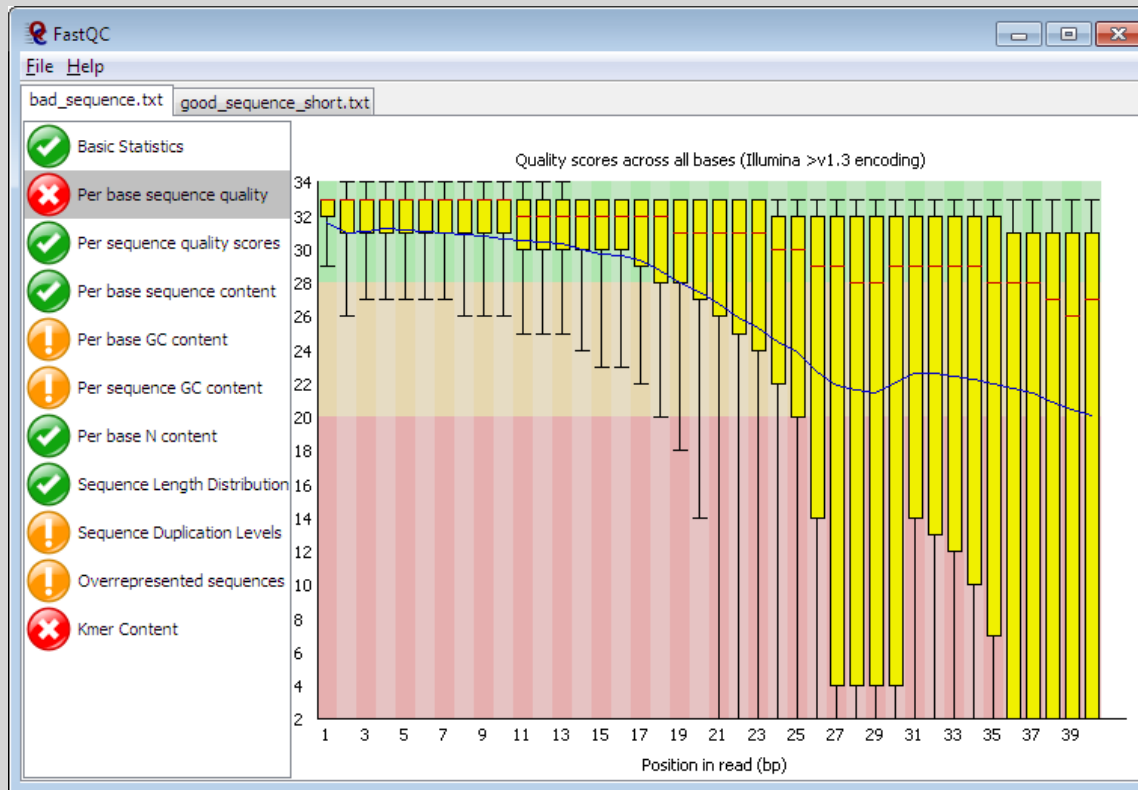
Per base sequence quality



FASTQC

FASTQC is one approach which provides a visual interpretation of the raw sequence reads

- <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>



Sequence Alignment

- Once sequence quality has been assessed, the next step is to align the sequence to a reference genome
- There are *many* distinct tools for doing this; which one you choose is often a reflection of your specific experiment and personal preference

BWA

Bowtie

SOAP2

Novoalign

mr/mrsFast

Eland

Blat

Bfast

BarraCUDA

CASHx

GSNAP

Mosiak

Stampy

SHRiMP

SeqMap

SLIDER

RMAP

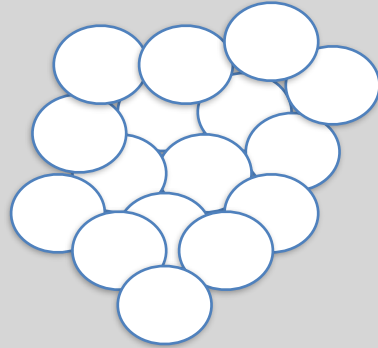
SSAHA

etc

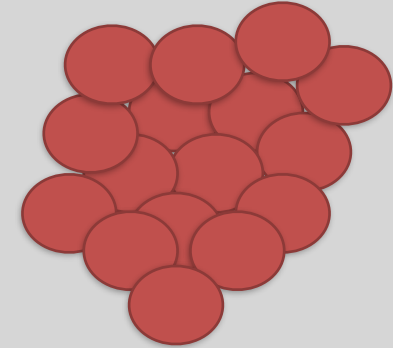
RNA Sequencing

The absolute basics

Normal Cells

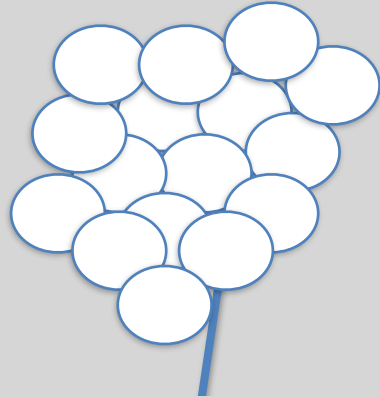


Mutated Cells

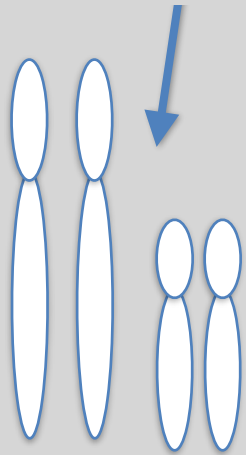


- The **mutated cells** behave differently than the **normal cells**
- We want to know what genetic mechanism is causing the difference
- One way to address this is to examine differences in gene expression via RNA sequencing...

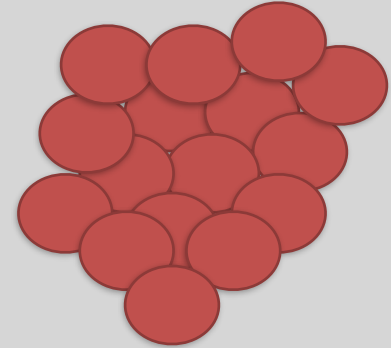
Normal Cells



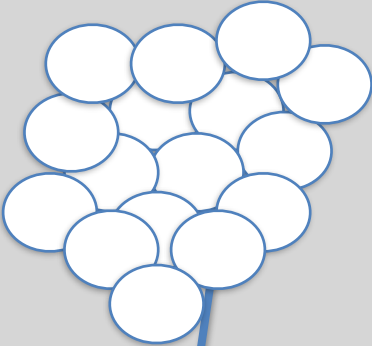
Each cell has a bunch of chromosomes



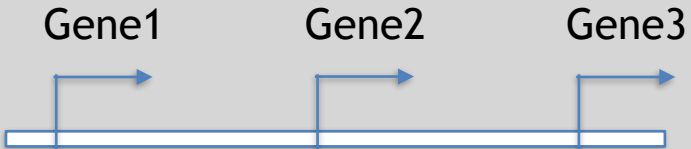
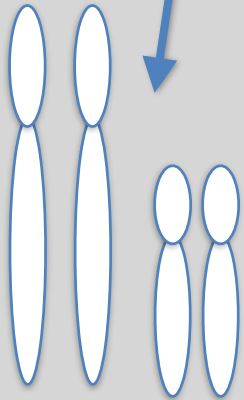
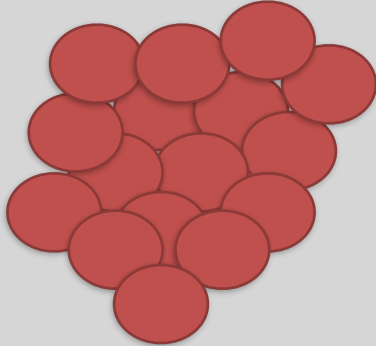
Mutated Cells



Normal Cells

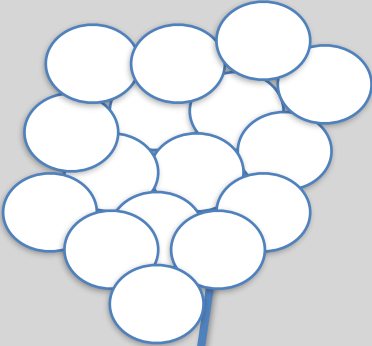


Mutated Cells

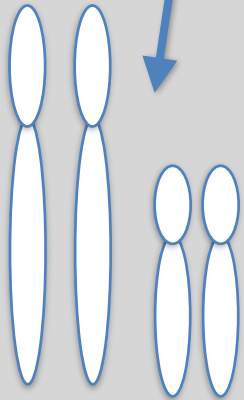
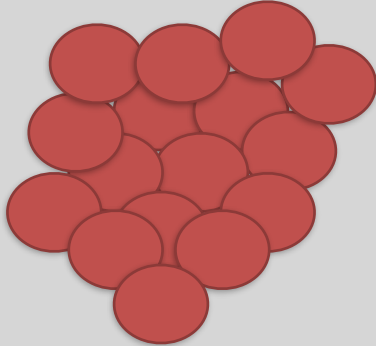


Each chromosome has a bunch of genes

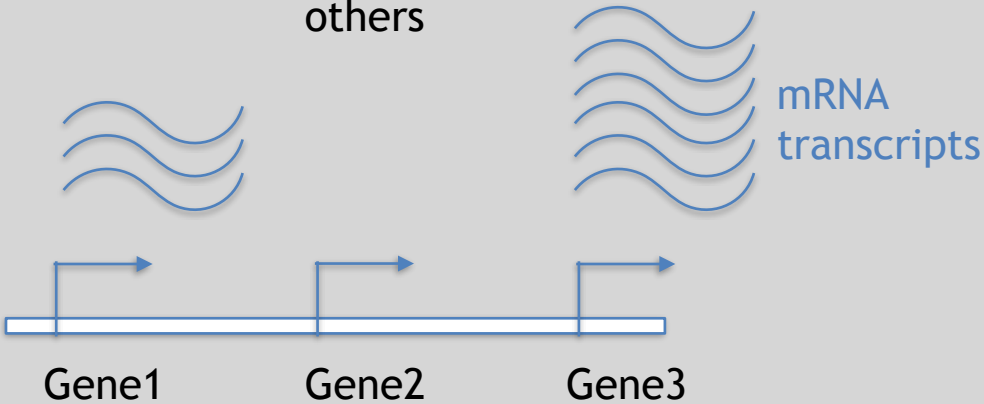
Normal Cells



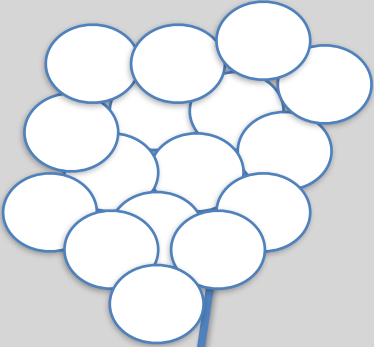
Mutated Cells



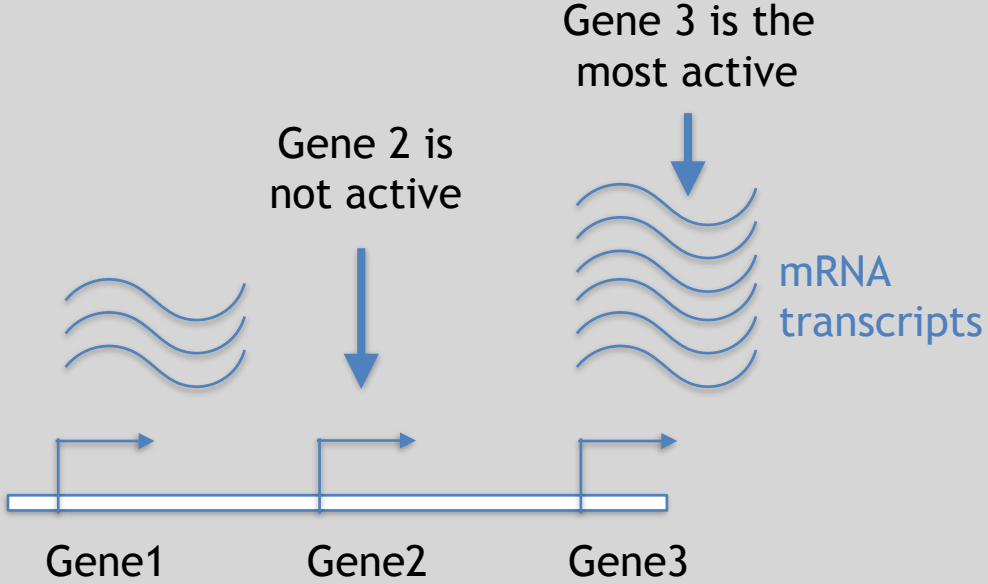
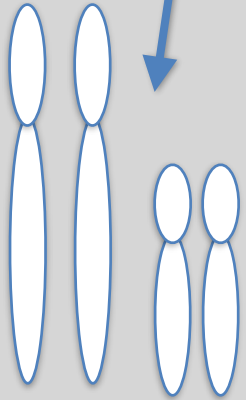
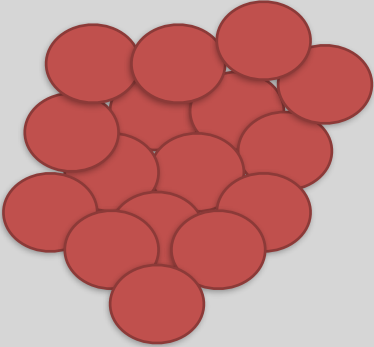
Some genes are active more than others



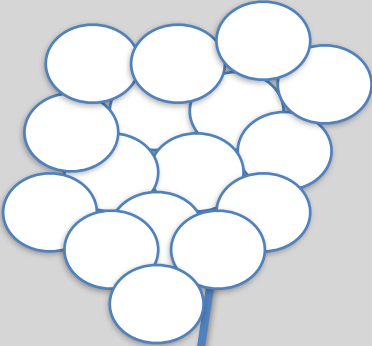
Normal Cells



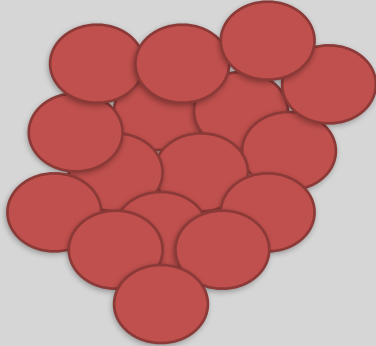
Mutated Cells



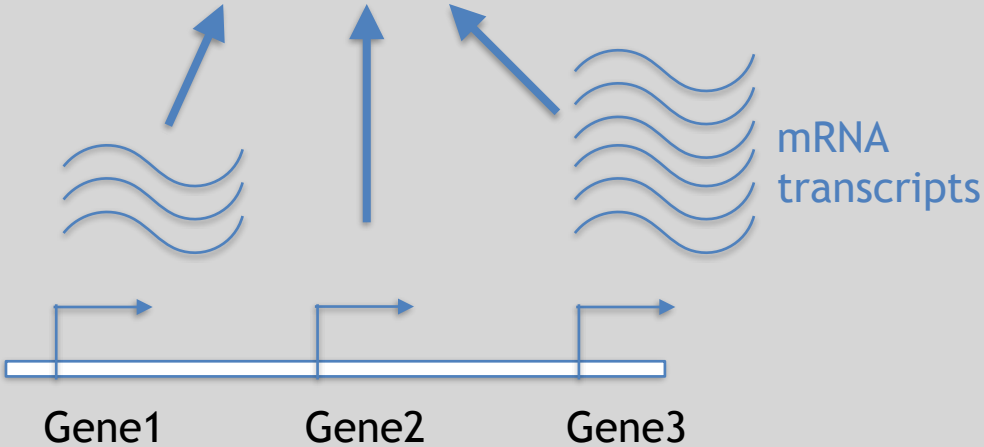
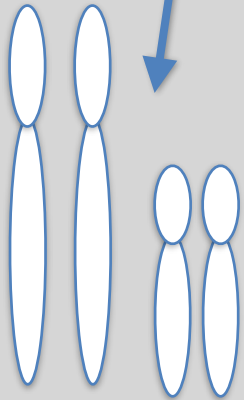
Normal Cells



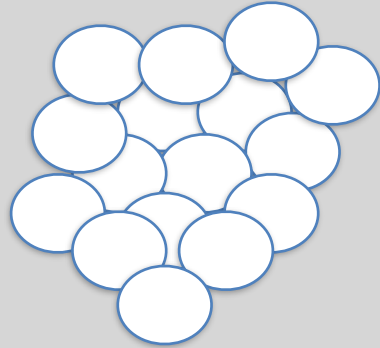
Mutated Cells



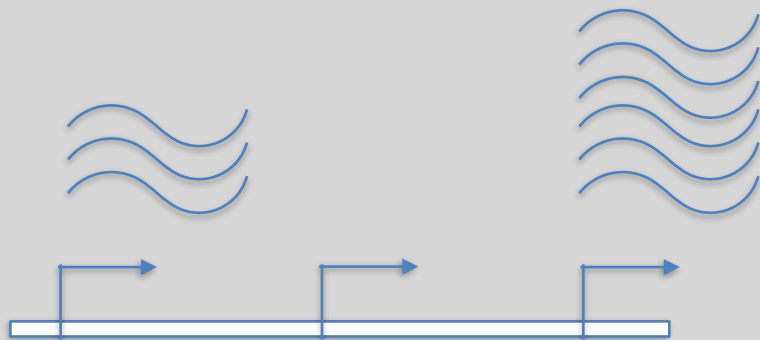
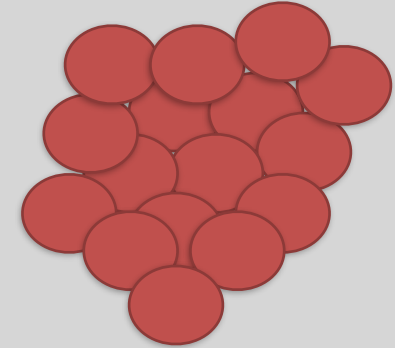
HTS tells us which genes are active, and how much they are transcribed!



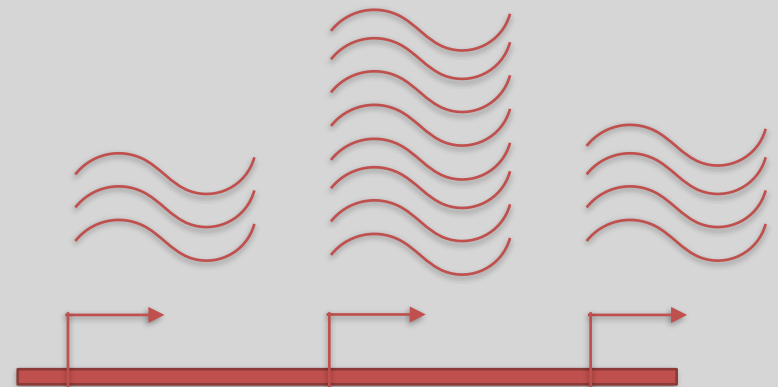
Normal Cells



Mutated Cells

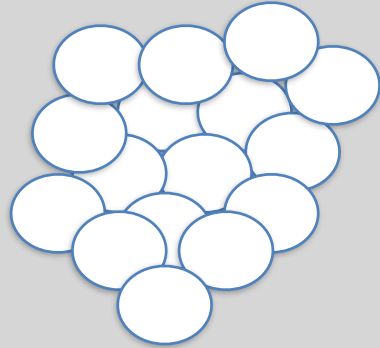


We use RNA-Seq to measure gene expression in normal cells ...

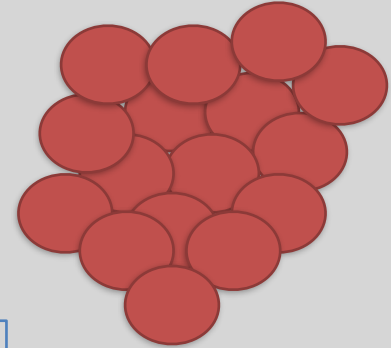


... then use it to measure gene expression in mutated cells

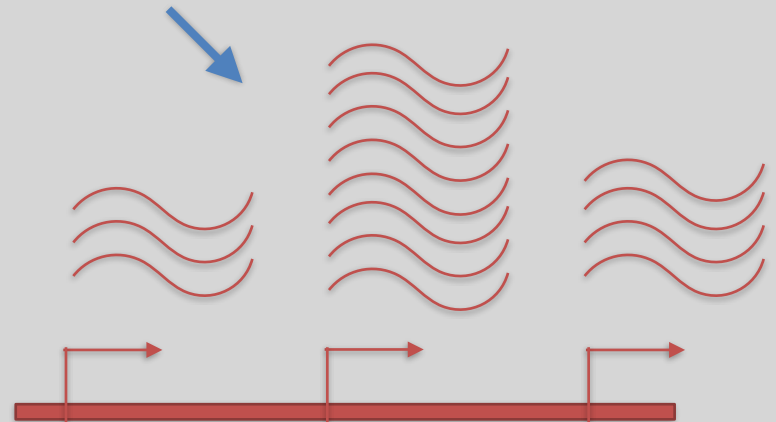
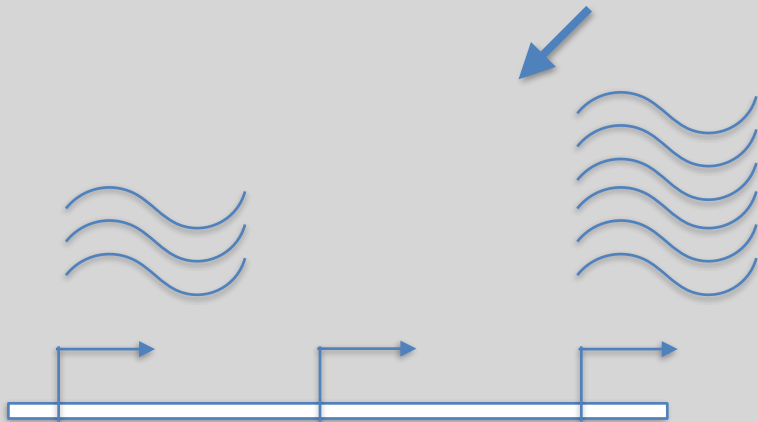
Normal Cells



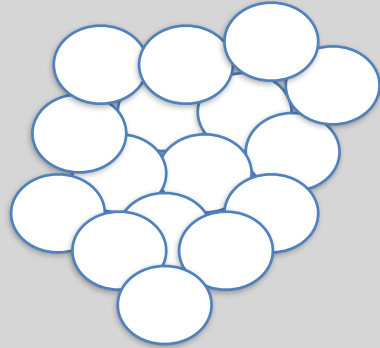
Mutated Cells



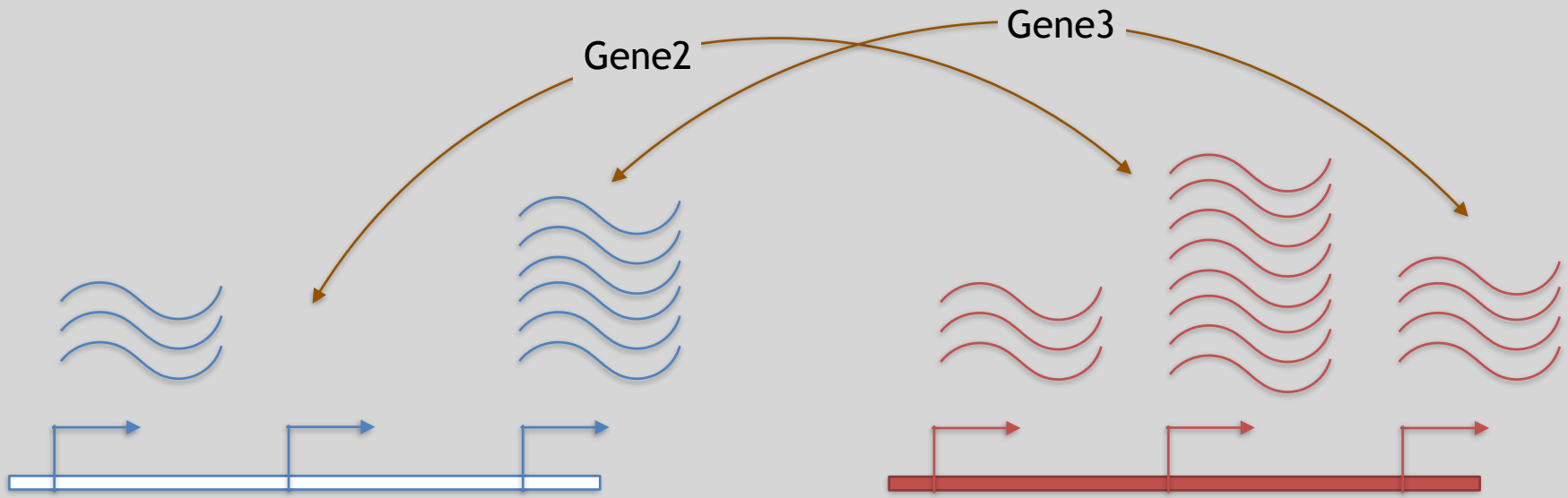
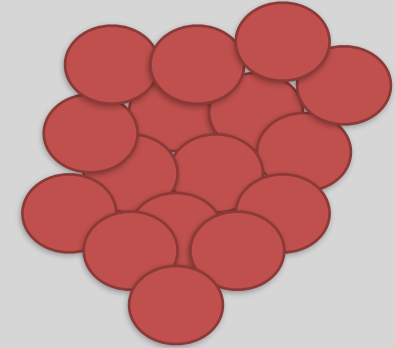
Then we can compare the two cell types to figure out what is different in the mutated cells!



Normal Cells



Mutated Cells



Differences apparent for Gene 2 and
to a lesser extent Gene 3

3 Main Steps for RNA-Seq:

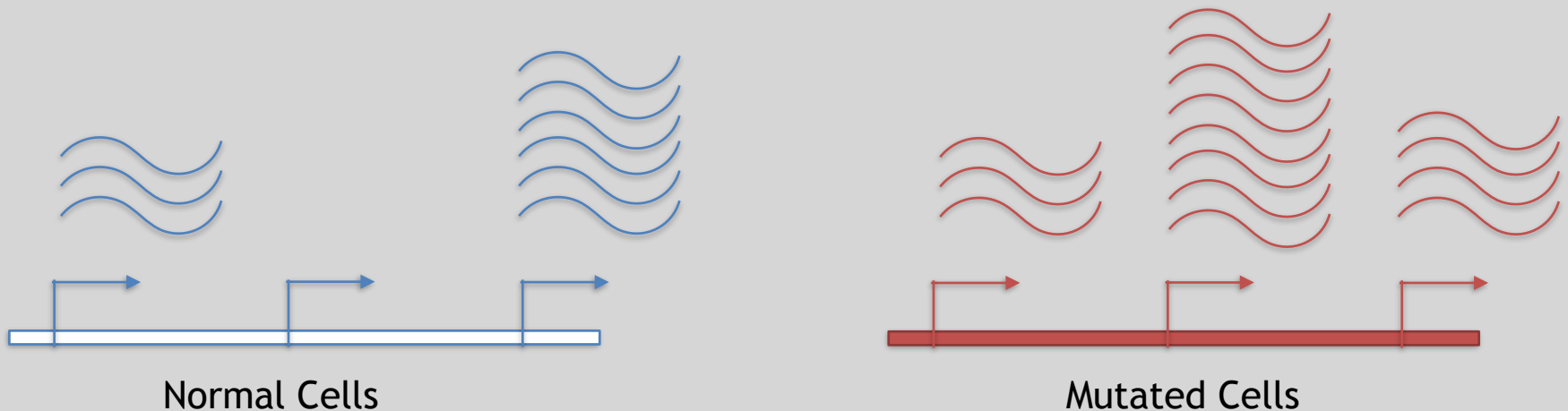
- 1) Prepare a sequencing library**
(RNA to cDNA conversion via reverse transcription)
- 2) Sequence**
(Using the same technologies as DNA sequencing)
- 3) Data analysis**
(Often the major bottleneck to overall success!)

We will discuss each of these steps in detail
(particularly the 3rd) next day!

Today we will get to the start of step 3!

Gene	WT-1	WT-2	WT-3	...
A1BG	30	5	13	...
AS1	24	10	18	...
...

We sequenced, aligned, counted the reads per gene in each sample to arrive at our data matrix



Do it Yourself!

Hand-on time!

https://bioboot.github.io/bimm143_F18/lectures/#14

Focus on **Sections 4** please
(After your Alignment is finished)

Feedback:

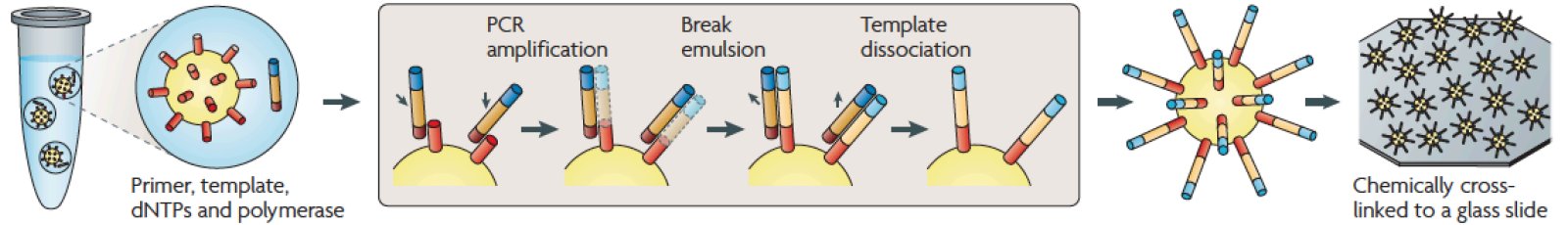
[Muddy Point Assessment]

Additional Reference Slides on SAM/BAM Format

Roche 454 - Pyrosequencing

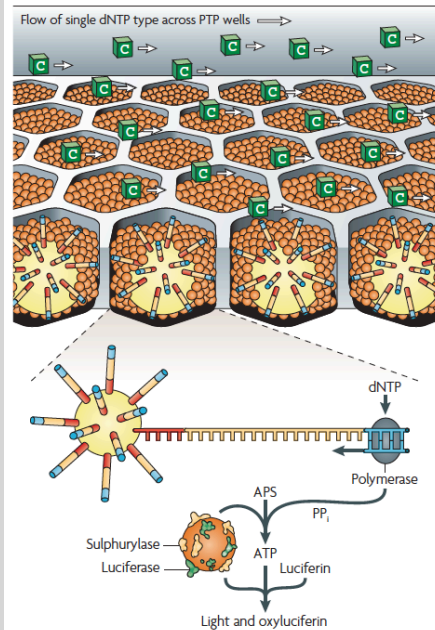
a Roche/454, Life/APG, Polonator Emulsion PCR

One DNA molecule per bead. Clonal amplification to thousands of copies occurs in microreactors in an emulsion



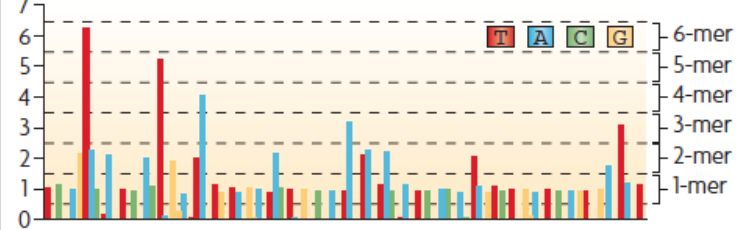
c Roche/454 — Pyrosequencing

1–2 million template beads loaded into PTP wells



d Flowgram

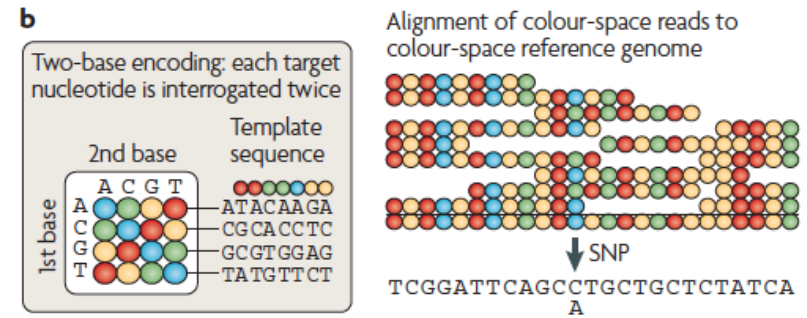
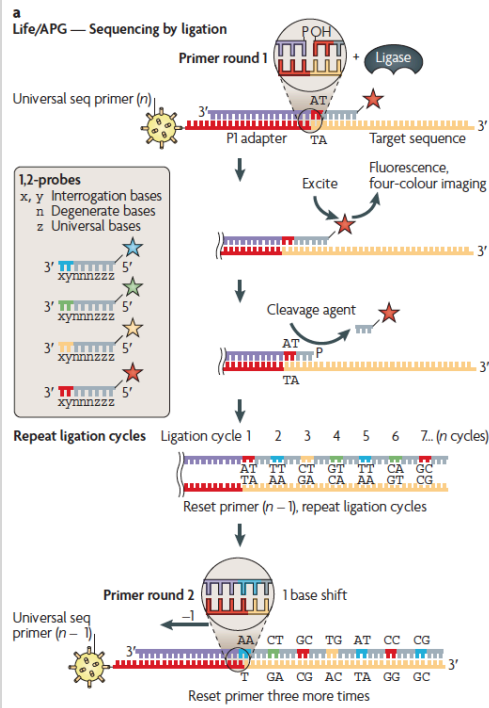
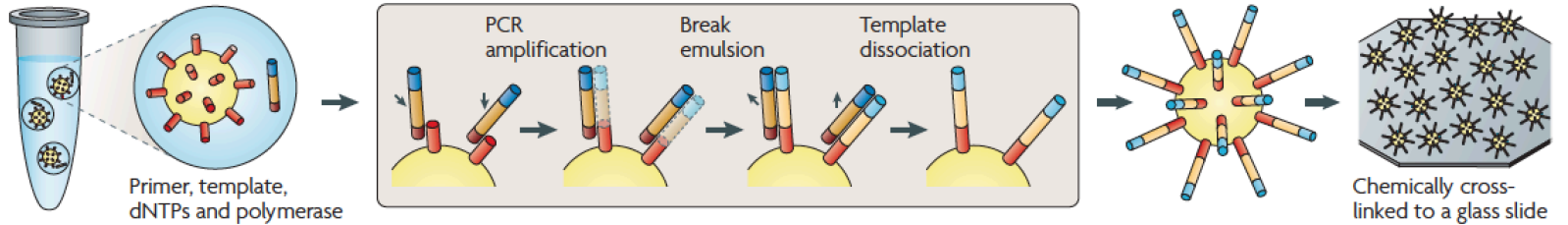
TCAGGTTTTTTTAAACAATCAACTTTTTGGATTAAAAATGTAGATAACTG
CATAAATTAATAACATCACATTAGTCTGATCAGTGAATTTAT



Life Technologies SOLiD - Sequence by Ligation

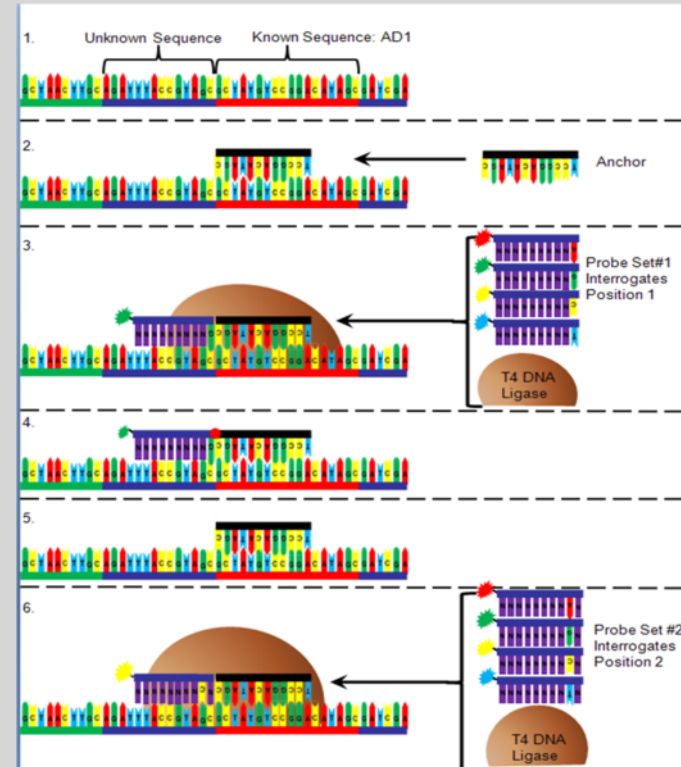
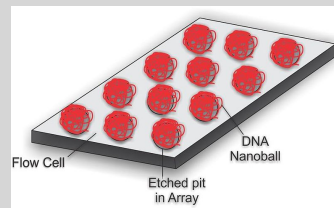
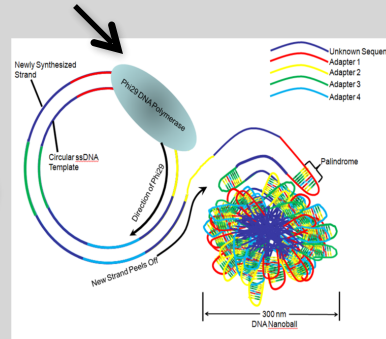
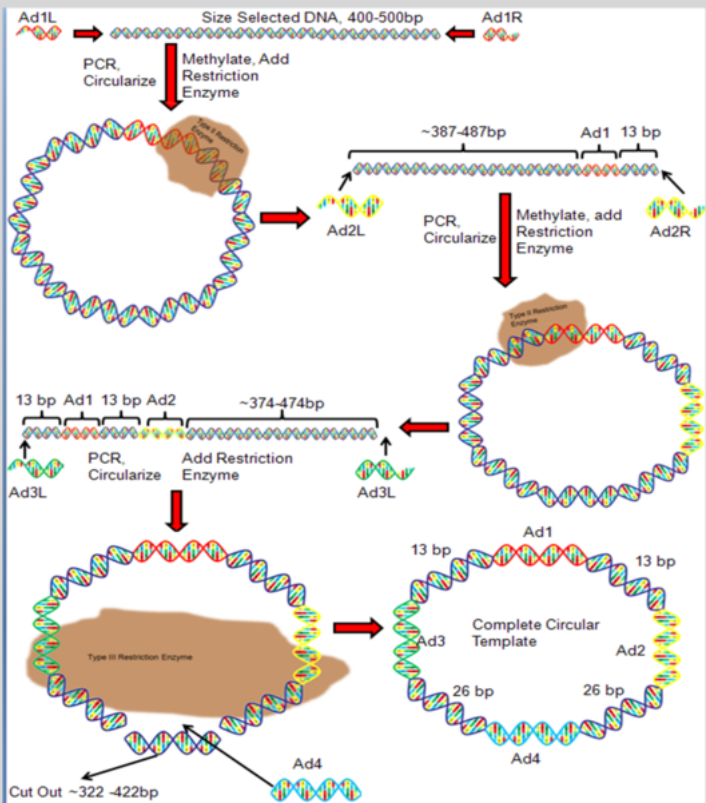
a Roche/454, Life/APG, Polonator Emulsion PCR

One DNA molecule per bead. Clonal amplification to thousands of copies occurs in microreactors in an emulsion



Complete Genomics - Nanoball Sequencing

Has proofreading ability!

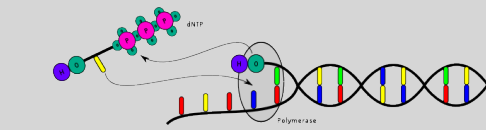


“Benchtop” Sequencers

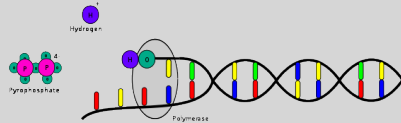
- Lower cost, lower throughput alternative for smaller scale projects
- Currently three significant platforms
 - Roche 454 GS Junior
 - Life Technology Ion Torrent
 - Personal Genome Machine (PGM)
 - Proton
 - Illumina MiSeq

Platform	List price	Approximate cost per run	Minimum throughput (read length)	Run time	Cost/Mb	Mb/h
454 GS Junior	\$108,000	\$1,100	35 Mb (400 bases)	8 h	\$31	4.4
Ion Torrent PGM						
(314 chip)	\$80,490 ^{a,b}	\$225 ^c	10 Mb (100 bases)	3 h	\$22.5	3.3
(316 chip)		\$425	100 Mb ^d (100 bases)	3 h	\$4.25	33.3
(318 chip)		\$625	1,000 Mb (100 bases)	3 h	\$0.63	333.3
MiSeq	\$125,000	\$750	1,500 Mb (2 × 150 bases)	27 h	\$0.5	55.5

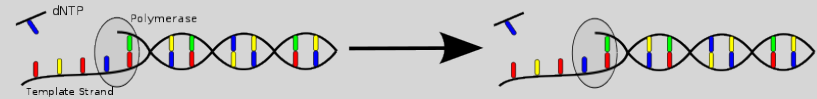
PGM - Ion Semiconductor Sequencing



Polymerase integrates a nucleotide.



Hydrogen and pyrophosphate are released.



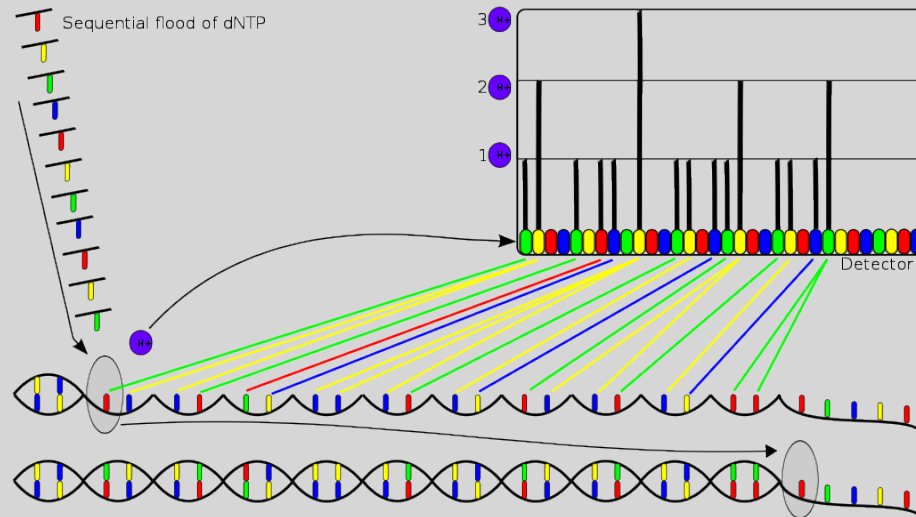
The nucleotide does not compliment the template - no release of hydrogen.



The nucleotide compliments the template - hydrogen is released.



The nucleotide compliments several bases in a row - multiple hydrogen ions are released.



Sequence Alignment

- Once sequence quality has been assessed, the next step is to align the sequence to a reference genome
- There are *many* distinct tools for doing this; which one you choose is often a reflection of your specific experiment and personal preference

BWA

Bowtie

SOAP2

Novoalign

mr/mrsFast

Eland

Blat

Bfast

BarraCUDA

CASHx

GSNAP

Mosiak

Stampy

SHRiMP

SeqMap

SLIDER

RMAP

SSAHA

etc

SAM Format

- Sequence Alignment/Map (SAM) format is the almost-universal sequence alignment format for NGS
 - binary version is BAM
- It consists of a header section (lines start with '@') and an alignment section
- The official specification can be found here:
 - <http://samtools.sourceforge.net/SAM1.pdf>

SAM header section

- Header lines contain vital metadata about the reference sequences, read and sample information, and (optionally) processing steps and comments.
- Each header line begins with an @, followed by a two-letter code that distinguishes the different type of metadata records in the header.
- Following this two-letter code are tab-delimited key-value pairs in the format **KEY:VALUE** (the SAM format specification names these tags and values).

https://bioboot.github.io/bimm143_F18/class-material/sam_format/

SAM Utilities

- **Samtools** is a common toolkit for analyzing and manipulating files in SAM/BAM format
 - <http://samtools.sourceforge.net/>
- **Picard** is a another set of utilities that can be used to manipulate and modify SAM files
 - <http://picard.sourceforge.net/>
- These can be used for viewing, parsing, sorting, and filtering SAM files as well as adding new information (e.g. Read Groups)

TODAYS MENU:

▶ **What is a Genome?**

- Genome sequencing and the Human genome project

▶ **What can we do with a Genome?**

- Comparative genomics

▶ **Modern Genome Sequencing**

- 1st, 2nd and 3rd generation sequencing

▶ **Workflow for NGS**

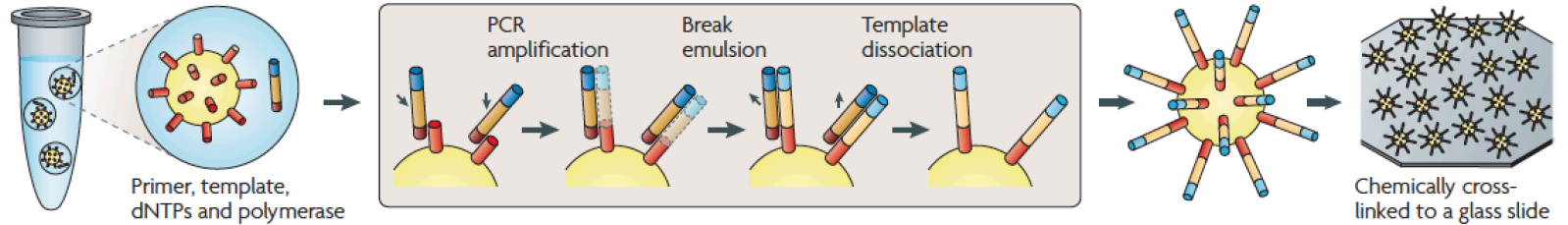
- RNA-Sequencing and discovering variation

Additional Reference Slides on Sequencing Methods

Roche 454 - Pyrosequencing

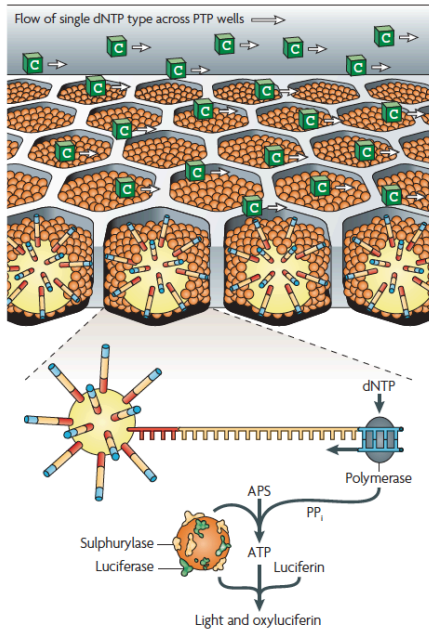
a Roche/454, Life/APG, Polonator Emulsion PCR

One DNA molecule per bead. Clonal amplification to thousands of copies occurs in microreactors in an emulsion



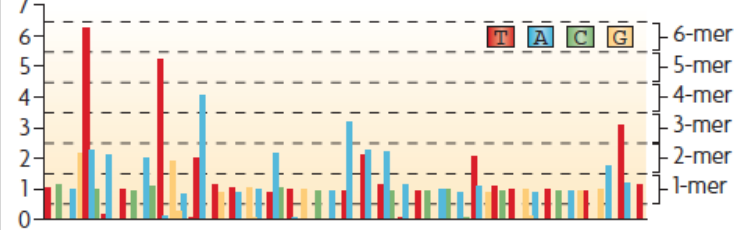
c Roche/454 — Pyrosequencing

1–2 million template beads loaded into PTP wells



d Flowgram

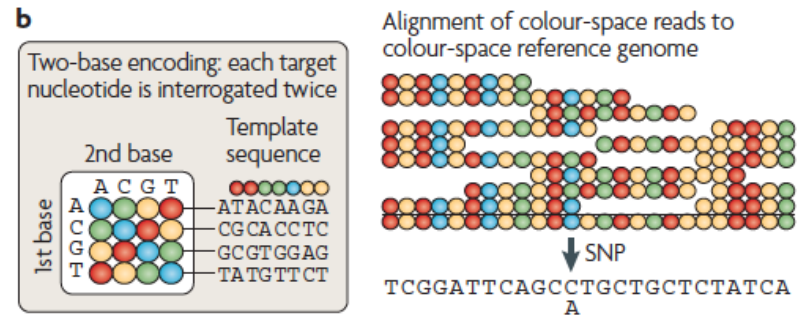
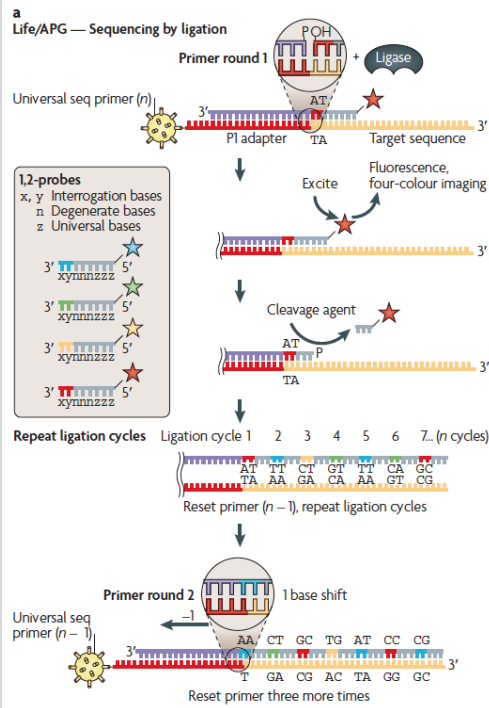
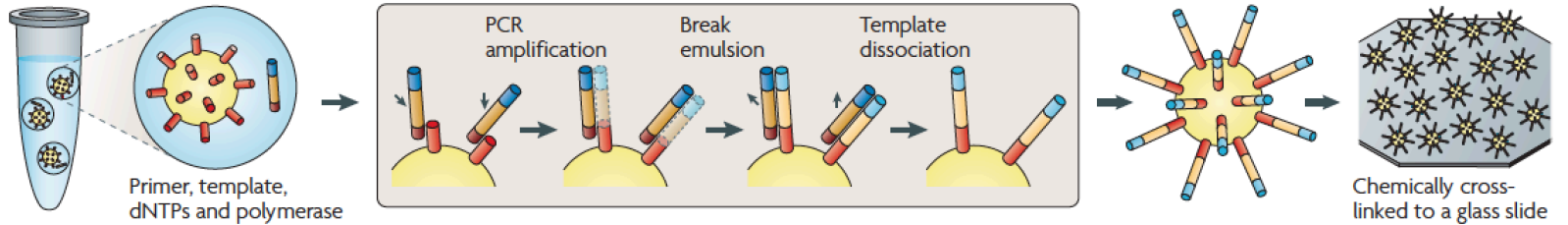
TCAGGTTTTTTTAAACAATCAACTTTTTGGATTAAAAATGTAGATAACTG
CATAAATTAATAACATCACATTAGTCTGATCAGTGAATTTAT



Life Technologies SOLiD - Sequence by Ligation

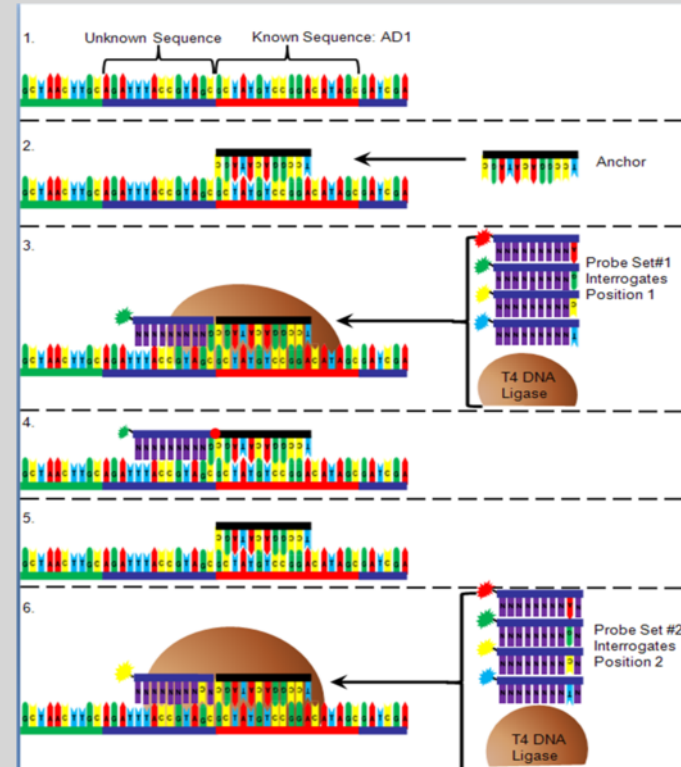
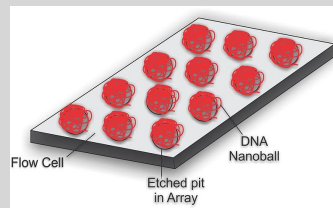
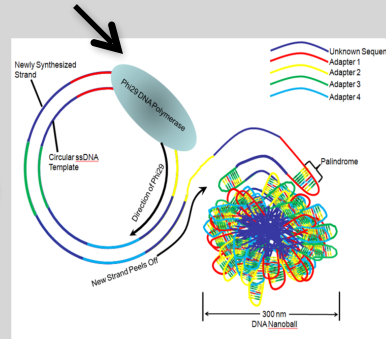
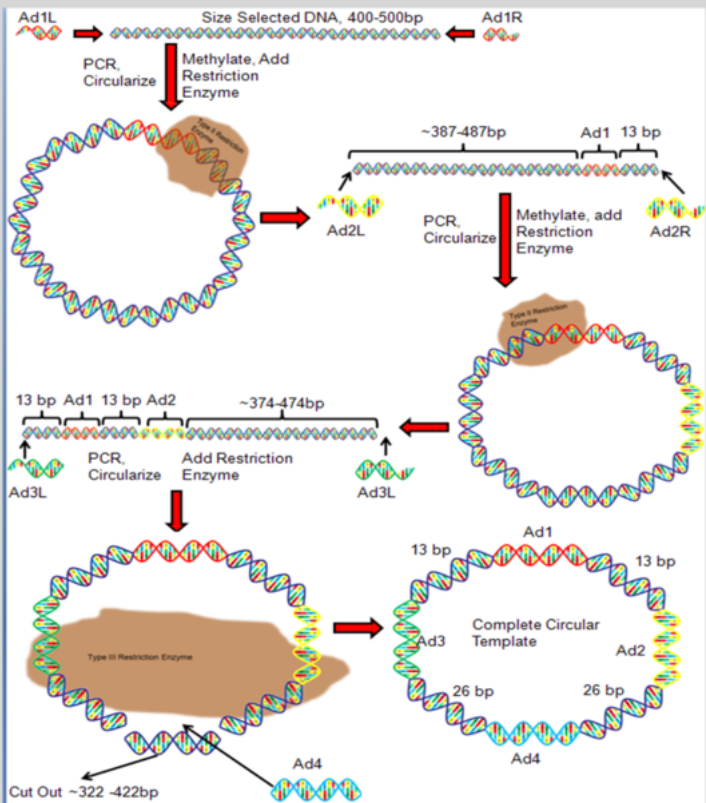
a Roche/454, Life/APG, Polonator Emulsion PCR

One DNA molecule per bead. Clonal amplification to thousands of copies occurs in microreactors in an emulsion



Complete Genomics - Nanoball Sequencing

Has proofreading ability!

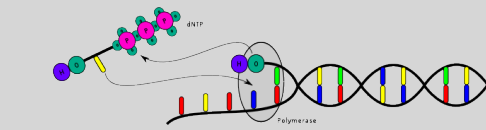


“Benchtop” Sequencers

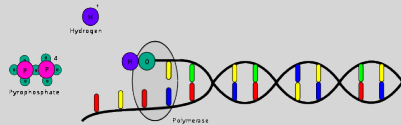
- Lower cost, lower throughput alternative for smaller scale projects
- Currently three significant platforms
 - Roche 454 GS Junior
 - Life Technology Ion Torrent
 - Personal Genome Machine (PGM)
 - Proton
 - Illumina MiSeq

Platform	List price	Approximate cost per run	Minimum throughput (read length)	Run time	Cost/Mb	Mb/h
454 GS Junior	\$108,000	\$1,100	35 Mb (400 bases)	8 h	\$31	4.4
Ion Torrent PGM						
(314 chip)	\$80,490 ^{a,b}	\$225 ^c	10 Mb (100 bases)	3 h	\$22.5	3.3
(316 chip)		\$425	100 Mb ^d (100 bases)	3 h	\$4.25	33.3
(318 chip)		\$625	1,000 Mb (100 bases)	3 h	\$0.63	333.3
MiSeq	\$125,000	\$750	1,500 Mb (2 × 150 bases)	27 h	\$0.5	55.5

PGM - Ion Semiconductor Sequencing



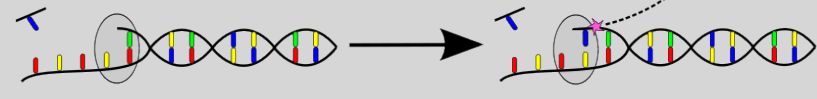
Polymerase integrates a nucleotide.



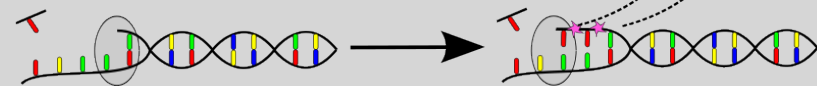
Hydrogen and pyrophosphate are released.



The nucleotide does not compliment the template - no release of hydrogen.



The nucleotide compliments the template - hydrogen is released.



The nucleotide compliments several bases in a row - multiple hydrogen ions are released.

