



# BGGN 213

## Genome Informatics I

Lecture 13

Barry Grant  
UC San Diego

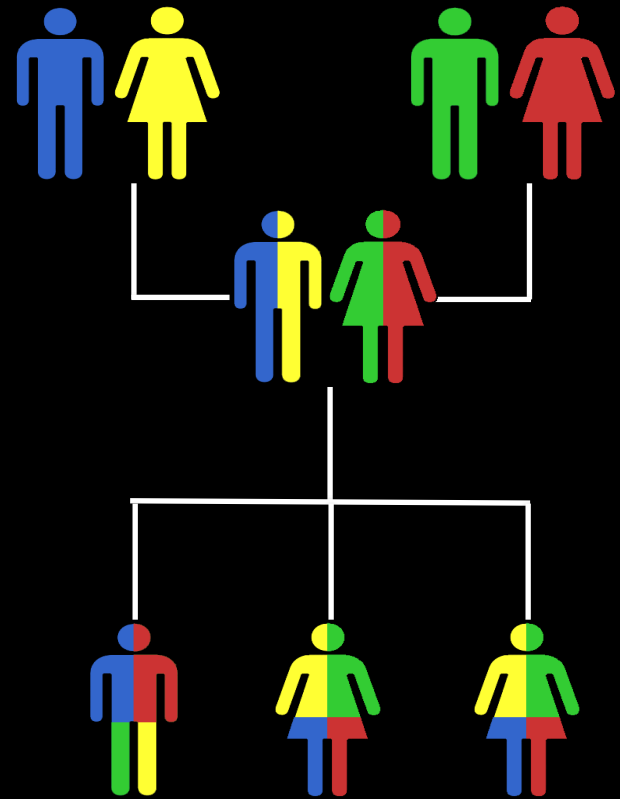
<http://thegrantlab.org/bggn213>

# Today's 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

# 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



# Genetics and Genomics

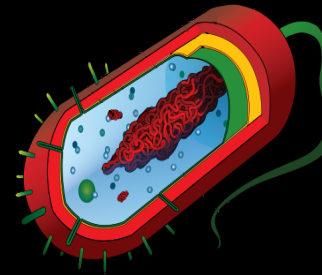
- **Genetics** is primarily the study of *individual genes*, mutations within those genes, and their inheritance patterns in order to understand specific traits.
- **Genomics** expands upon classical genetics and considers aspects of the *entire genome*, typically using computer aided approaches.



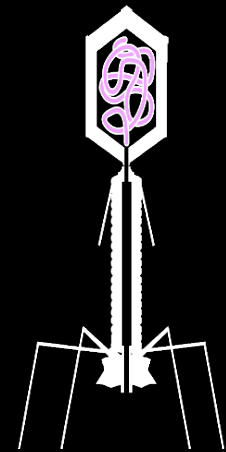
# Genomes come in many shapes

Side note!

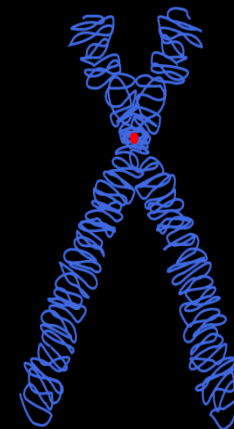
- Primarily DNA, but can be RNA in the case of some viruses
- Some genomes are circular, others linear
- Can be organized into discrete units (chromosomes) or freestanding molecules (plasmids)



Prokaryote



Bacteriophage



Eukaryote

Side note!

## CHROMOSOMES CLOSE-UP

Chromosomes consist largely of double-helical DNA. Cells package the DNA into the nucleus by wrapping it around “spools” composed of histone proteins. The DNA-protein combination is known as chromatin. (Each color represents one chromosome.)

Under a microscope, a Eukaryotic cell's genome (i.e. collection of chromosomes) resembles a chaotic jumble of noodles. The looping is not random however and appears to play a role in controlling gene regulation.

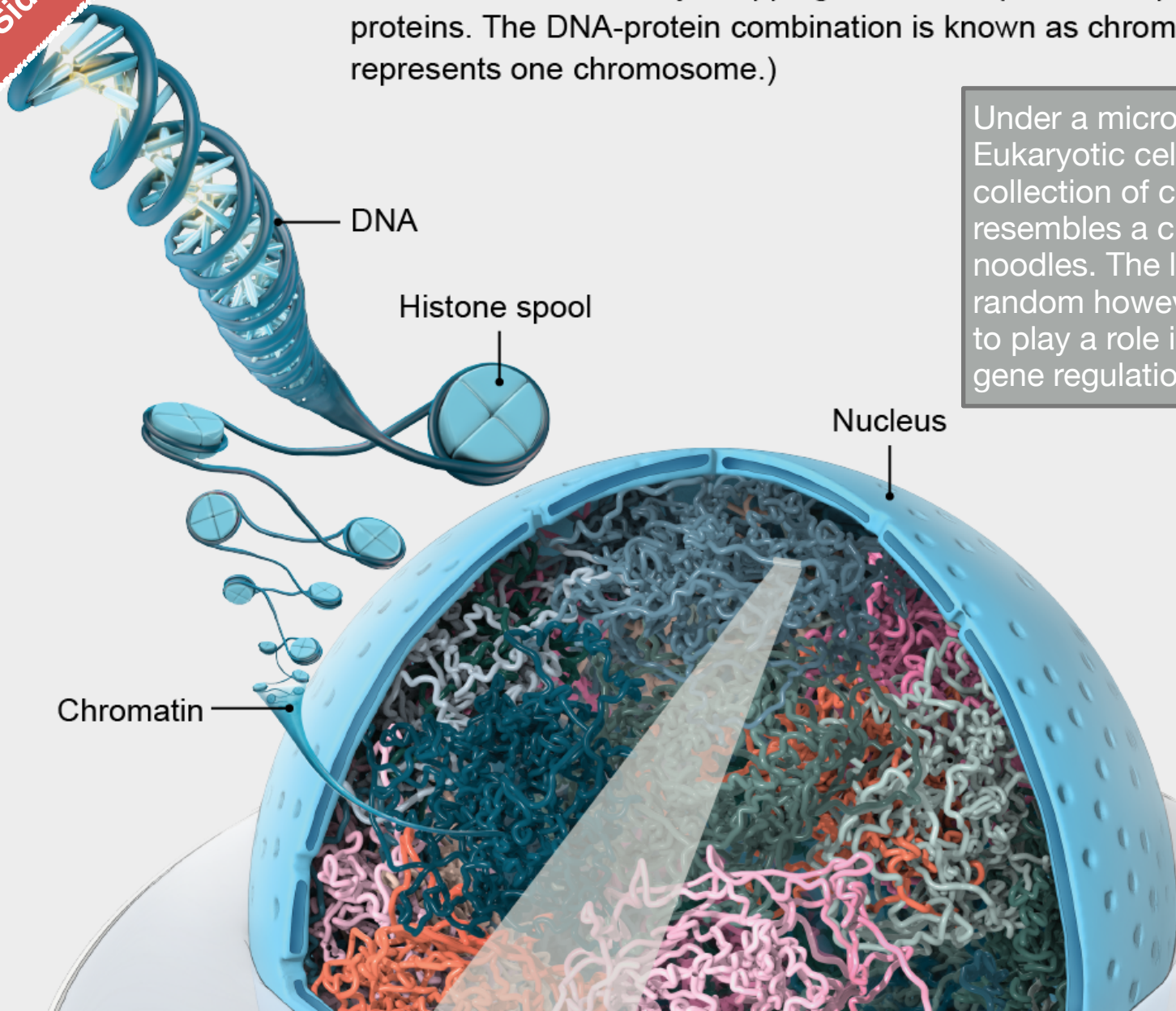
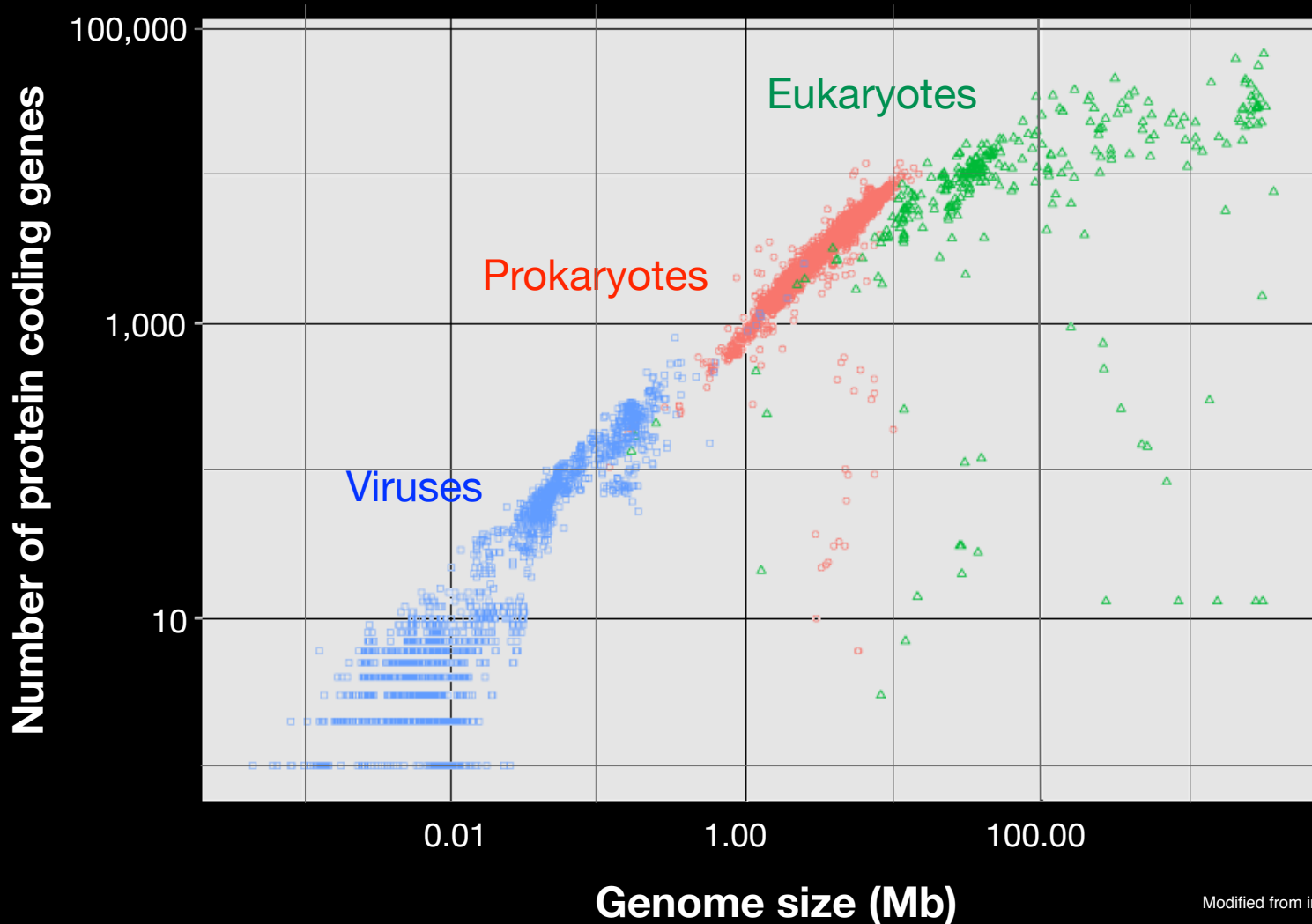


Image credit:  
[Scientific American](#)  
[March 2019](#)

# 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

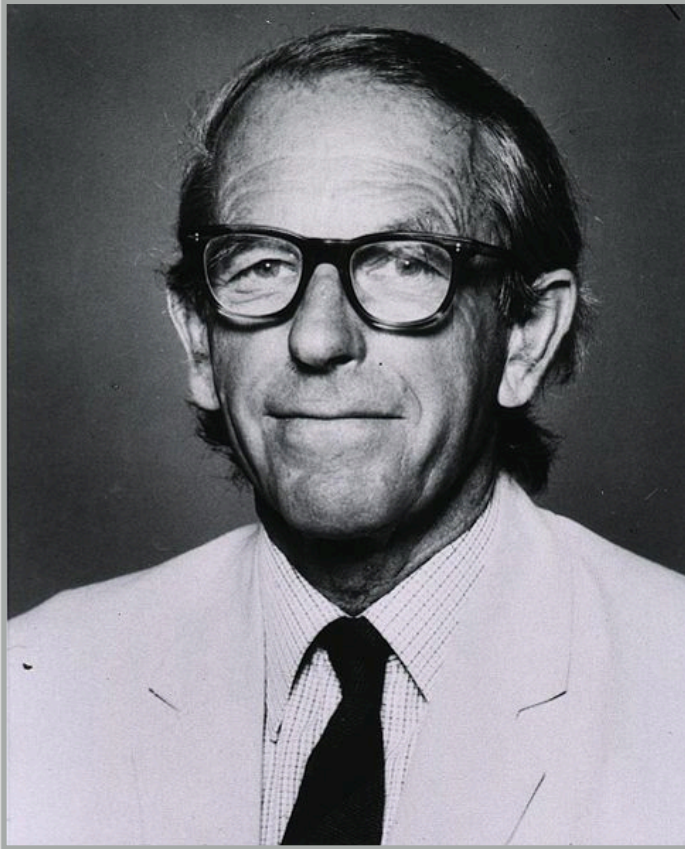
This resource organizes information on genomes including sequences, maps, chromosomes, assemblies, and annotations.

<b>Using Genome</b> <ul style="list-style-type: none"><li><a href="#">Help</a></li><li><a href="#">Browse by Organism</a></li><li><a href="#">Download / FTP</a></li><li><a href="#">Download FAQ</a></li><li><a href="#">Submit a genome</a></li></ul>	<b>Custom resources</b> <ul style="list-style-type: none"><li><a href="#">Human Genome</a></li><li><a href="#">Microbes</a></li><li><a href="#">Organelles</a></li><li><a href="#">Viruses</a></li><li><a href="#">Prokaryotic reference genomes</a></li></ul>	<b>Other Resources</b> <ul style="list-style-type: none"><li><a href="#">Assembly</a></li><li><a href="#">BioProject</a></li><li><a href="#">BioSample</a></li><li><a href="#">Map Viewer</a></li><li><a href="#">Protein Clusters</a></li></ul>
<b>Genome Tools</b> <ul style="list-style-type: none"><li><a href="#">BLAST the Human Genome</a></li><li><a href="#">Microbial Nucleotide BLAST</a></li><li><a href="#">TaxPlot (3-way Genome Comparison)</a></li></ul>	<b>Genome Annotation and Analysis</b> <ul style="list-style-type: none"><li><a href="#">Eukaryotic Genome Annotation</a></li><li><a href="#">Prokaryotic Genome Annotation</a></li><li><a href="#">PASC (Pairwise Sequence Comparison)</a></li></ul>	<b>External Resources</b> <ul style="list-style-type: none"><li><a href="#">GOLD - Genomes Online Database</a></li><li><a href="#">Ensembl Genome Browser</a></li><li><a href="#">Bacteria Genomes at Sanger</a></li><li><a href="#">Large-Scale Genome Sequencing (NHGRI)</a></li></ul>

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

<b>GETTING STARTED</b> <ul style="list-style-type: none"><li>NCBI Education</li><li>NCBI Help Manual</li><li>NCBI Handbook</li><li>Training &amp; Tutorials</li></ul>	<b>RESOURCES</b> <ul style="list-style-type: none"><li>Chemicals &amp; Bioassays</li><li>Data &amp; Software</li><li>DNA &amp; RNA</li><li>Domains &amp; Structures</li><li>Genes &amp; Expression</li><li>Genetics &amp; Medicine</li><li>Genomes &amp; Maps</li><li>Homology</li><li>Literature</li><li>Proteins</li><li>Sequence Analysis</li><li>Taxonomy</li><li>Training &amp; Tutorials</li><li>Variation</li></ul>	<b>POPULAR</b> <ul style="list-style-type: none"><li>PubMed</li><li>Bookshelf</li><li>PubMed Central</li><li>PubMed Health</li><li>BLAST</li><li>Nucleotide</li><li>Genome</li><li>SNP</li><li>Gene</li><li>Protein</li><li>PubChem</li></ul>	<b>FEATURED</b> <ul style="list-style-type: none"><li>Genetic Testing Registry</li><li>PubMed Health</li><li>GenBank</li><li>Reference Sequences</li><li>Gene Expression Omnibus</li><li>Map Viewer</li><li>Human Genome</li><li>Mouse Genome</li><li>Influenza Virus</li><li>Primer-BLAST</li><li>Sequence Read Archive</li></ul>	<b>NCBI INFORMATION</b> <ul style="list-style-type: none"><li>About NCBI</li><li>Research at NCBI</li><li>NCBI News</li><li>NCBI FTP Site</li><li>NCBI on Facebook</li><li>NCBI on Twitter</li><li>NCBI on YouTube</li></ul>
---	--	---	--	--

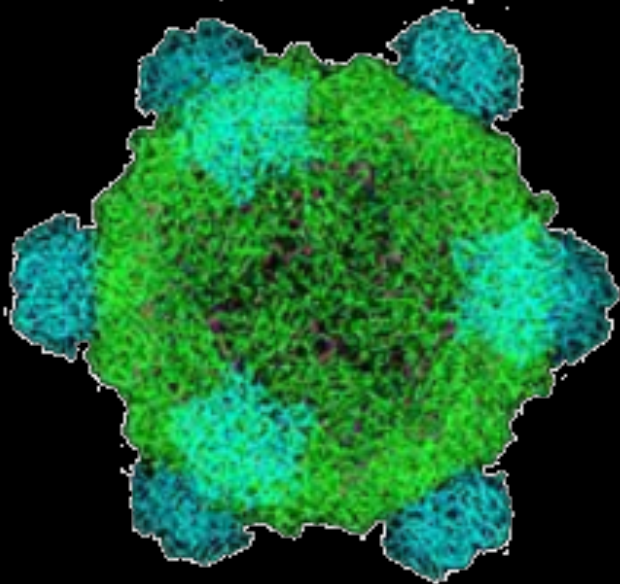
# Early Genome Sequencing



- 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

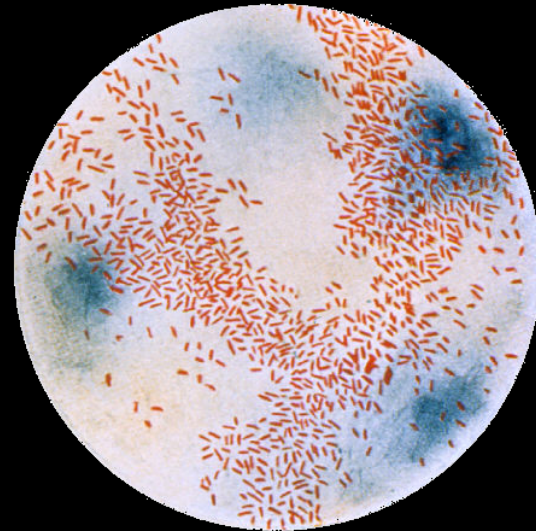


# The First Sequenced Genomes



## Bacteriophage $\phi$ -X174

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



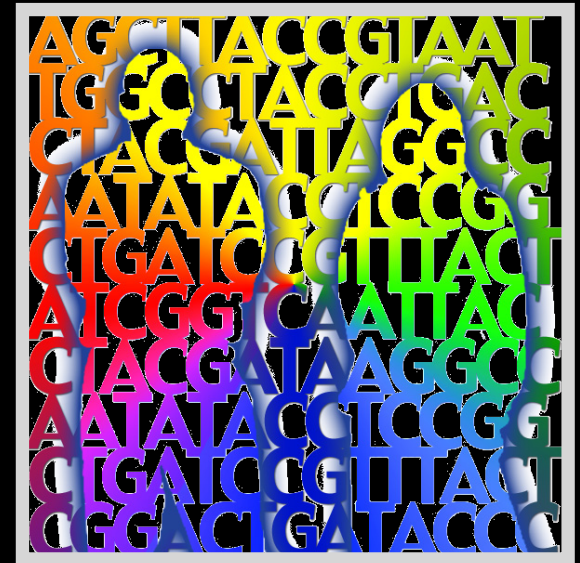
## Haemophilus influenzae

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



# The Human Genome Project

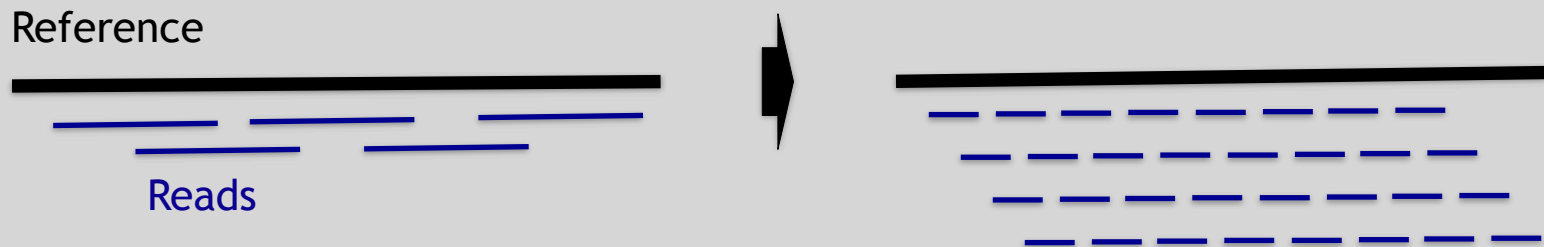
- The Human Genome Project (HGP) was an international, public consortium that began in 1990
  - Initiated by James Watson
  - Primarily led by Francis Collins
  - Eventual Cost: \$2.7 Billion
- Celera Genomics was a private corporation that started in 1998
  - Headed by Craig Venter
  - Eventual Cost: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



HHMI

# Modern Genome Sequencing

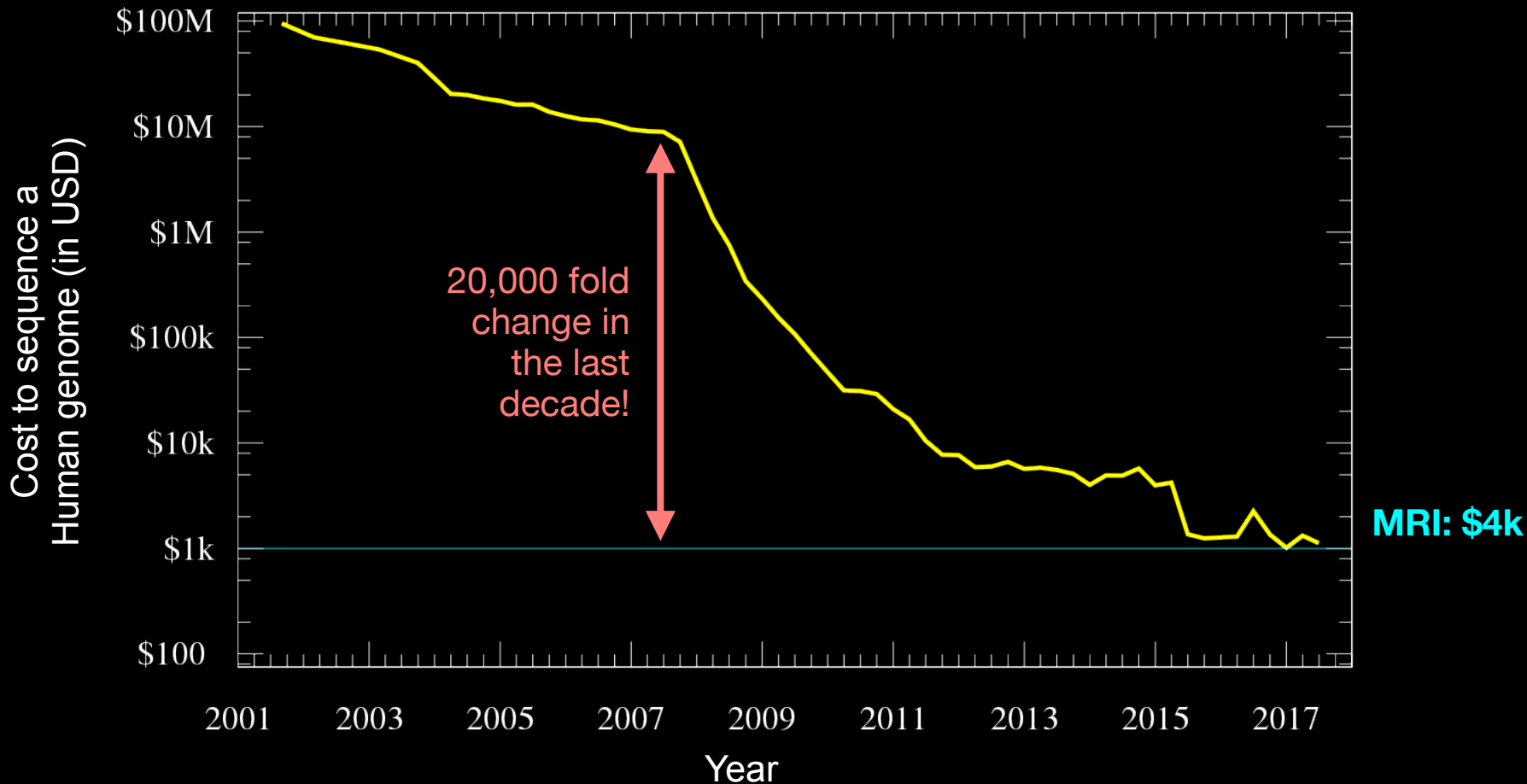
- Next Generation Sequencing (NGS) technologies have resulted in a paradigm shift from long reads at low coverage to short reads at high coverage
- This provides numerous opportunities for new and expanded genomic applications



# Rapid progress of genome sequencing



# Rapid progress of genome sequencing



# 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)

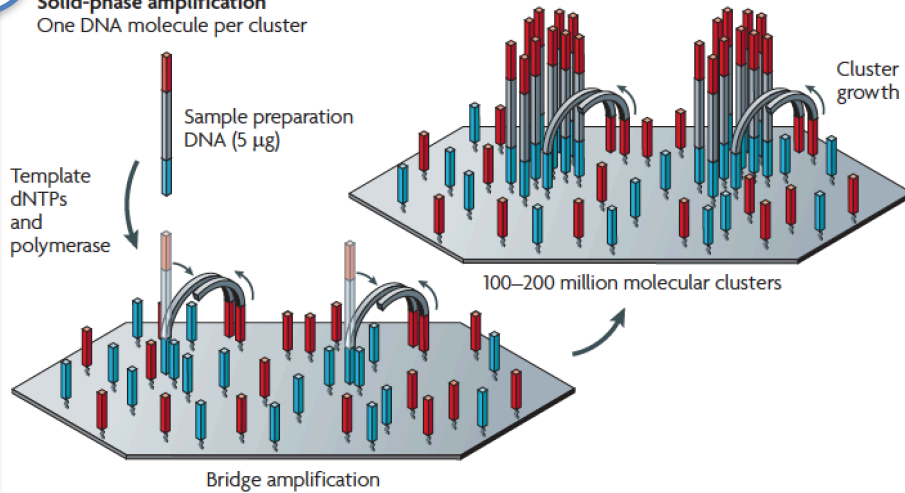
# Modern NGS Sequencing Platforms

	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

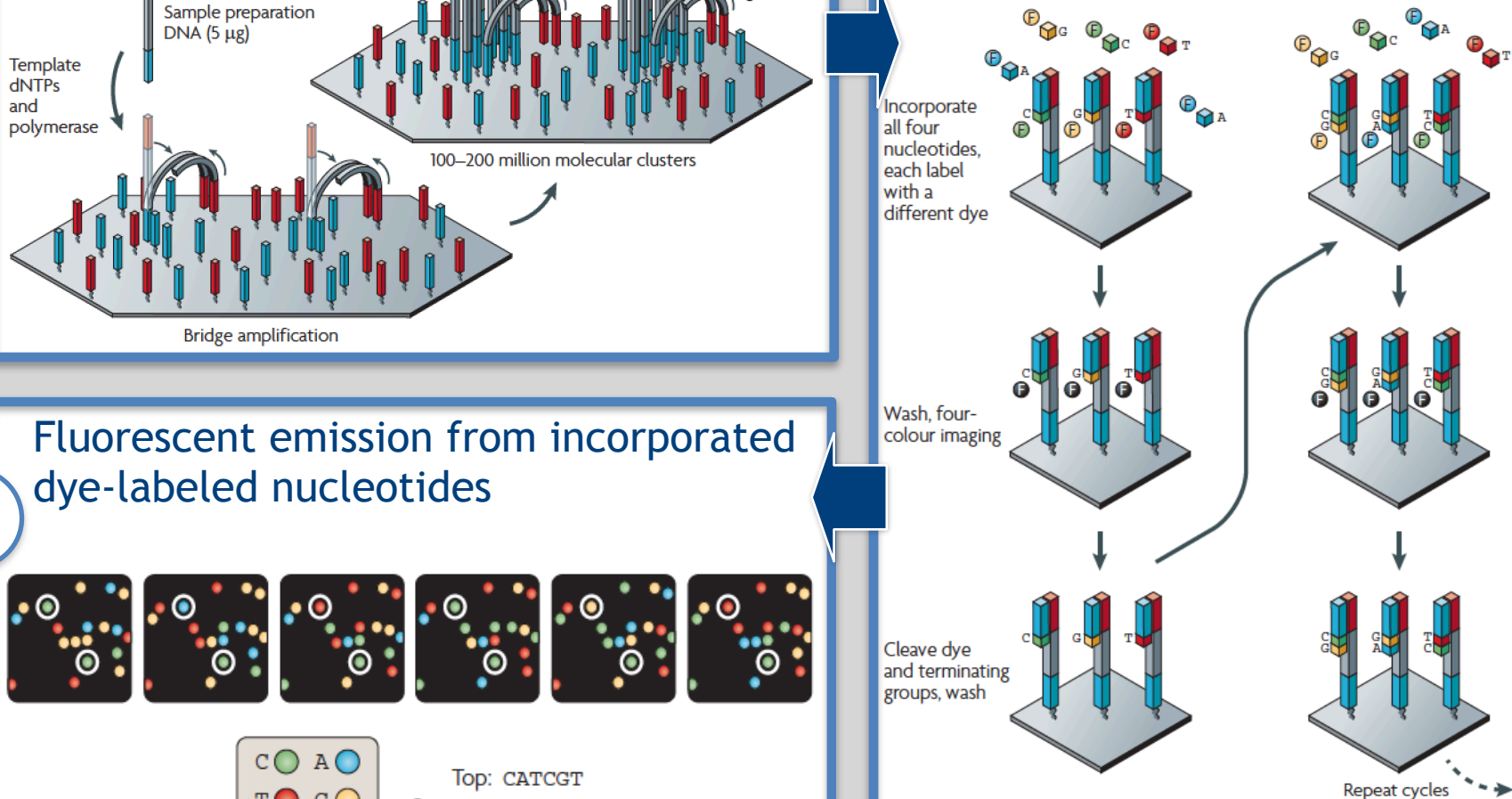
Illumina/Solexa — Reversible terminators

Incorporate all four nucleotides, each label with a different dye

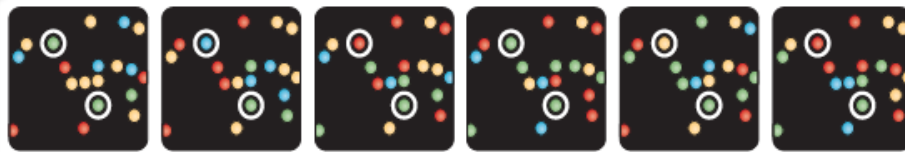
Wash, four-colour imaging

Cleave dye and terminating groups, wash

Cleave dye & blocking group, repeat...

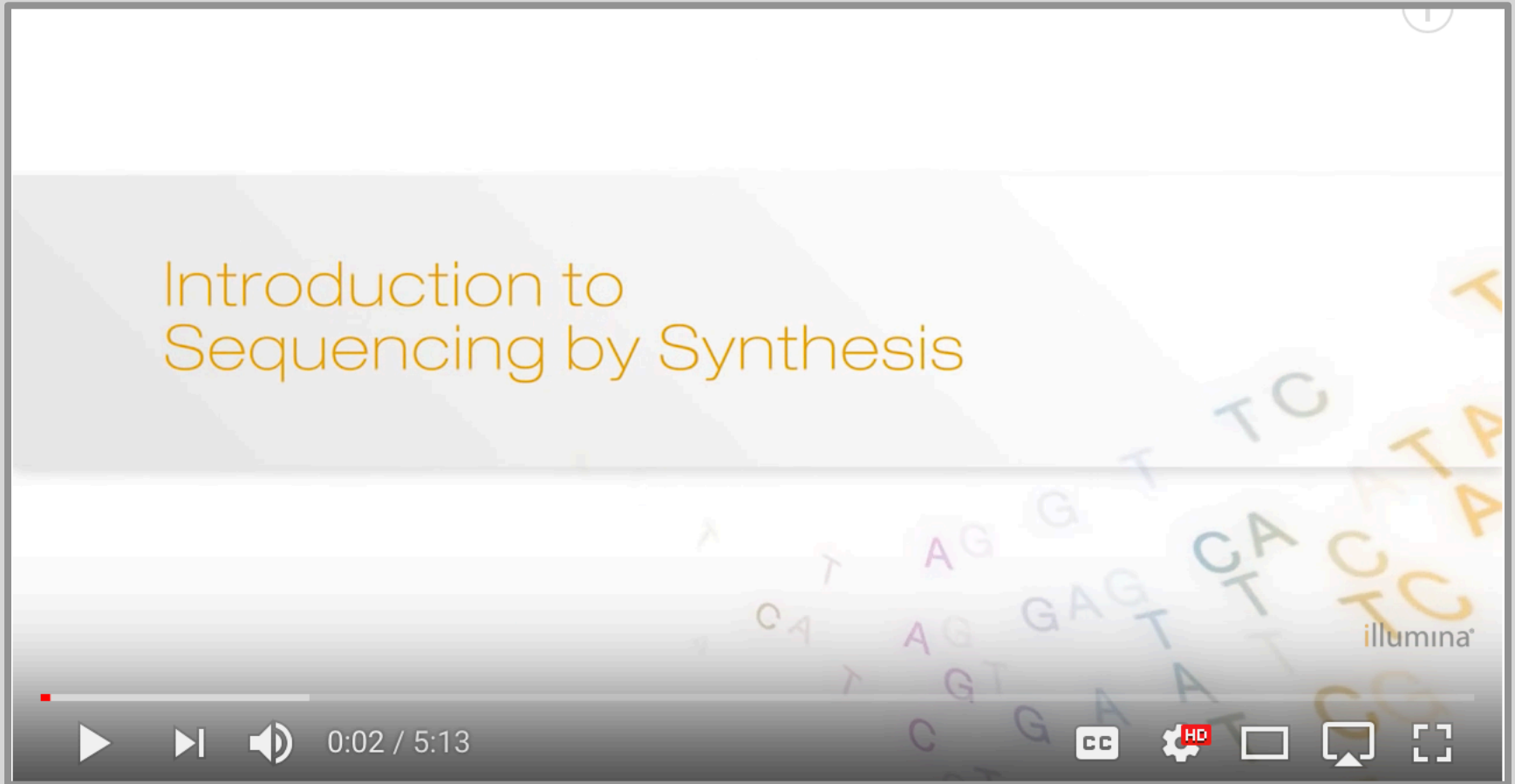


## 3 Fluorescent emission from incorporated dye-labeled nucleotides



Top: CATCGT  
Bottom: CCCCCC

# Illumina Sequencing - Video

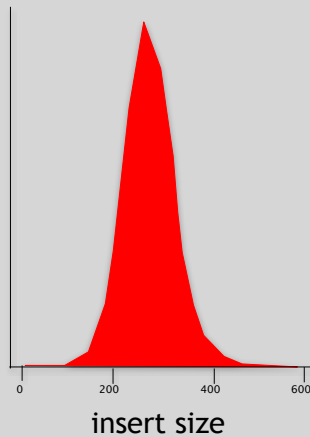
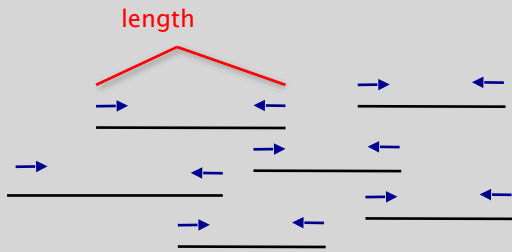


[https://www.youtube.com/watch?src\\_vid=womKfikWlxM&v=fCd6B5HRaZ8](https://www.youtube.com/watch?src_vid=womKfikWlxM&v=fCd6B5HRaZ8)

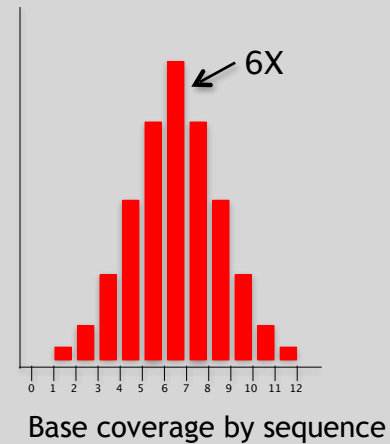
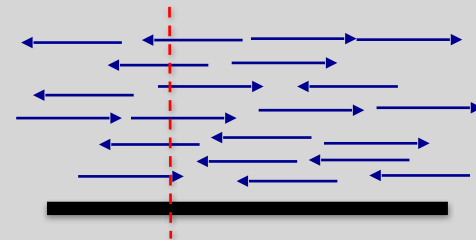


# NGS Sequencing Terminology

Insert Size



Sequence Coverage



# Summary: “Generations” of DNA Sequencing

	First generation	Second generation <sup>a</sup>	Third generation <sup>a</sup>
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 “3<sup>rd</sup>” 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 & BioStars

Side-Note:

A good repository of analysis software can be found at <http://seqanswers.com> and <https://www.biostars.org/>

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

## 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
<a href="#">4peaks</a>	Allows viewing sequencing trace files, motif searching trimming, BLAST and exporting sequences.	<a href="#">Sequencing</a>	Sequence analysis			Freeware	Mac OS X
<a href="#">AB Large Indel Tool</a>	Identifies deviations in clone insert size that indicate intra-chromosomal structural variations compared to a reference genome.	<a href="#">InDel discovery</a> <a href="#">Sequencing</a>	Mapping		Perl	GPL	Linux 64
<a href="#">AB Small Indel Tool</a>	The SOLID™ Small Indel Tool processes the indel evidences found in the pairing step of the SOLID™ System Analysis Pipeline Tool (Corona Lite).	<a href="#">InDel discovery</a> <a href="#">Sequencing</a>	Mapping Alignment		Perl C++	GPL	Linux 64
<a href="#">ABBA</a>	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	<a href="#">Genomic Assembly</a>	Assembly Scaffolding			Artistic License	Linux
<a href="#">ABMapper</a>	Maps RNA-Seq reads to target genome considering possible multiple mapping locations and splice junctions	<a href="#">Genomics</a> <a href="#">Transcriptomics</a>	Mapping Alignment		C++ Perl	GPLv3	Linux
<a href="#">ABySS</a>	ABySS is a de novo sequence assembler designed for short reads and large genomes.	<a href="#">De-novo assembly</a>	Assembly De Bruijn graph	MPI OpenMP	C++	Free for academic use	POSIX Linux Mac OS X
<a href="#">Ardanter Removal</a>	Removes ardanter fragments from raw short read	<a href="#">General</a>	Ardanter Removal	Trimming	Java	Custom Licence	Linux 64

Log In

Wiki navigation: [Main page](#), [Recent changes](#), [Random page](#), [Help](#)

Software: [Software hub](#), [Browse software](#), [Software list](#)

Toolbox: [What links here](#), [Related changes](#), [Special pages](#), [Printable version](#), [Permanent link](#), [Browse properties](#)

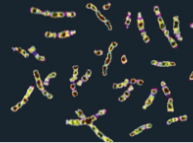
**What can we do with all  
this sequence information?**

# 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”

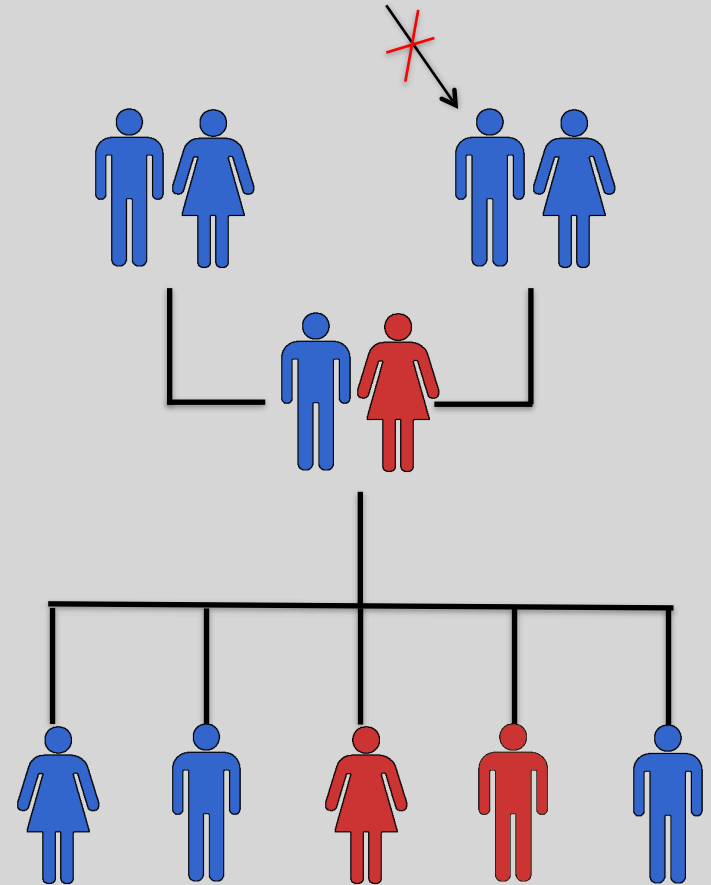
-Kruglyak & Nickerson, 2001

- 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.

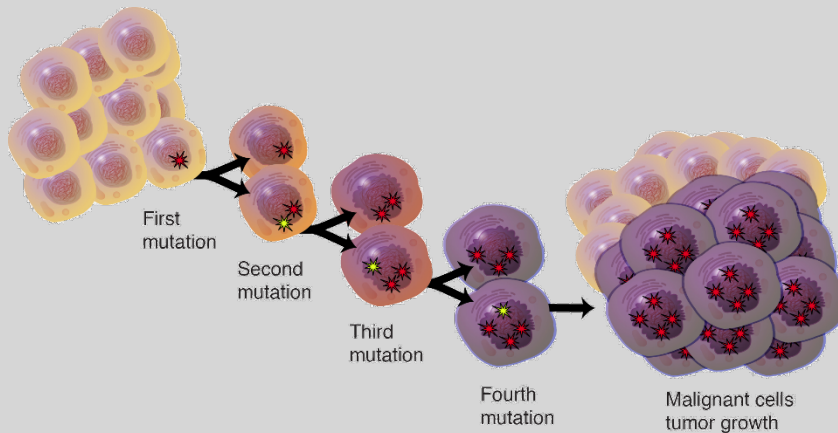


# Germline Variation

- 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



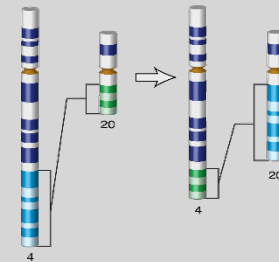
- Mutations in non-germline cells that are not passed along to offspring
- Can occur during mitosis or from the environment itself
- Are an integral part in tumor progression and evolution

# 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

AATCTGAGGCAT  
AATCTCAGGCAT

AATCTGAGGCAT  
AATCT--AGGCAT



# 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!

# Discovering Variation: SNPs and INDELS

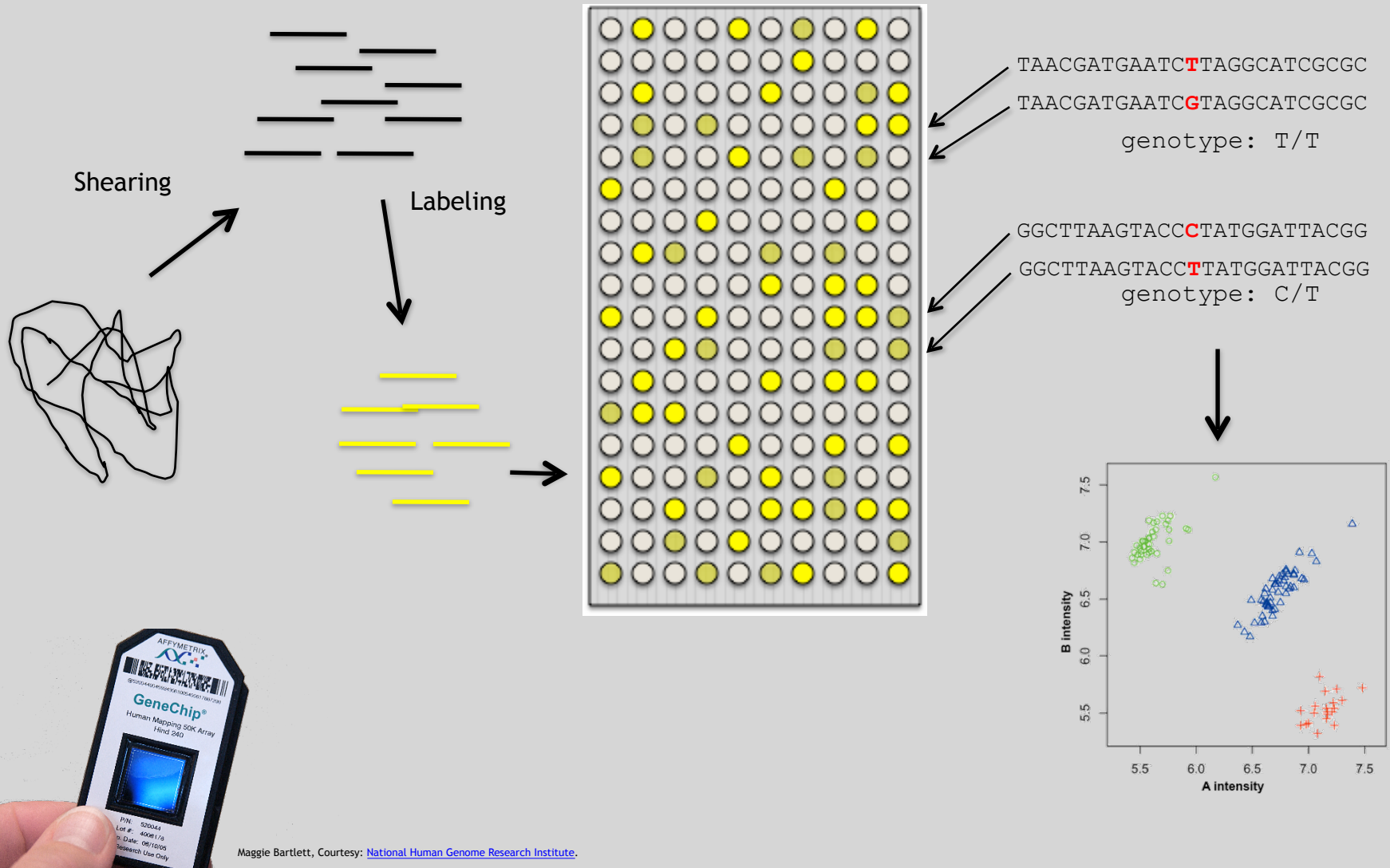
- Small variants require the use of sequence data to initially be discovered
- Most approaches align sequences to a reference genome to identify differing positions
- The amount of DNA sequenced is proportional to the number of times a region is covered by a sequence read
  - More sequence coverage equates to more support for a candidate variant site



# 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

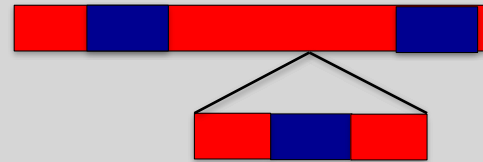
There are numerous ways genetic variation can exhibit functional effects

Premature stop codons



TAC -> TAA

Gene or exon deletion



Frameshift mutation



TAC -> T-C

Oct-1

Transcription factor binding disruption



ATGCAAAT -> ATGCAGAT

Do it Yourself!

# Hand-on time!

[https://bioboot.github.io/bgggn213\\_S19/lectures/#13](https://bioboot.github.io/bgggn213_S19/lectures/#13)

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, showing three levels of genomic detail for a region on Chromosome 12 (GRCh38.p3).

**Chromosome image:** Shows the whole chromosome with a highlighted region of interest (12:25,204,789-25,250,936) and haplotypes and patches.

**Overview image:** Shows a detailed view of the region, including chromosome bands, contigs, and genes. A callout highlights a "Gene or region of interest" (RPI1-29515.4). The Gene Legend indicates merged Ensembl/Havana, processed transcript, RNA gene, and pseudogene.

**Zoomable Region image:** Shows a highly detailed view of the region, including transcripts (splice variants) and genes. A callout highlights "Transcripts (splice variants)" for KRAS-002 and KRAS-003.

**Configuration and Navigation:** The left sidebar contains options for "Location-based displays", "Genetic Variation", "Markers", and "Other genes". Callouts indicate "Change or add data tracks" for the "Configure this page" and "Bookmark this page" options. The "Location" field is set to 12:25204789-25250936, and the "Gene" field is empty.

# 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

Meqablast 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

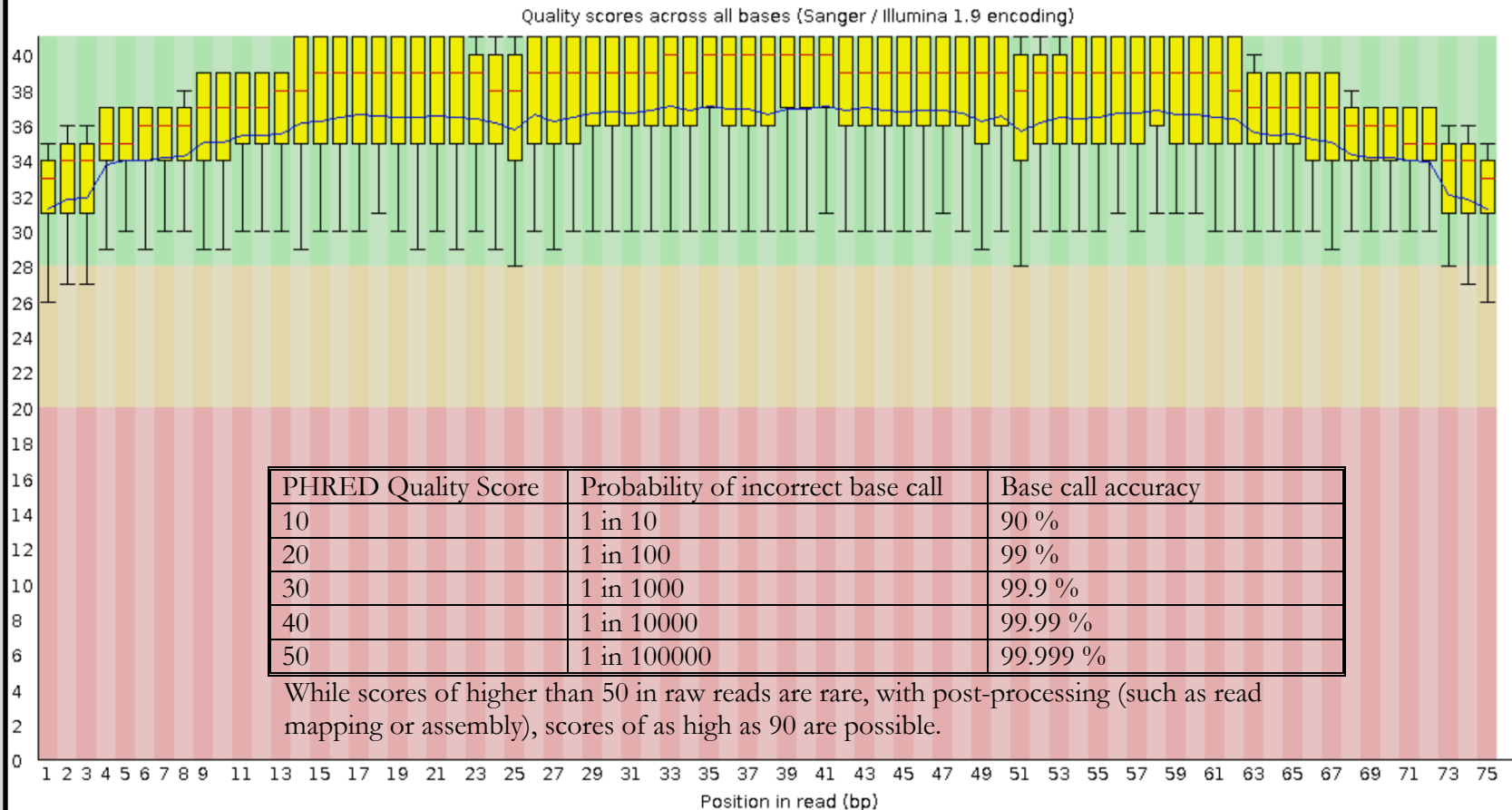
		ASCII Range	Offset	Score Range
Sanger, Illumina (Ver > 1.8)	fastqsanger	33-126	33	0-93
Solexa, Illumina (Ver < 1.3)	fastqsolexa	59-126	64	5-62
Illumina (Ver 1.3 -1.7)	fastqillumina	64-126	64	0-62

```
> 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

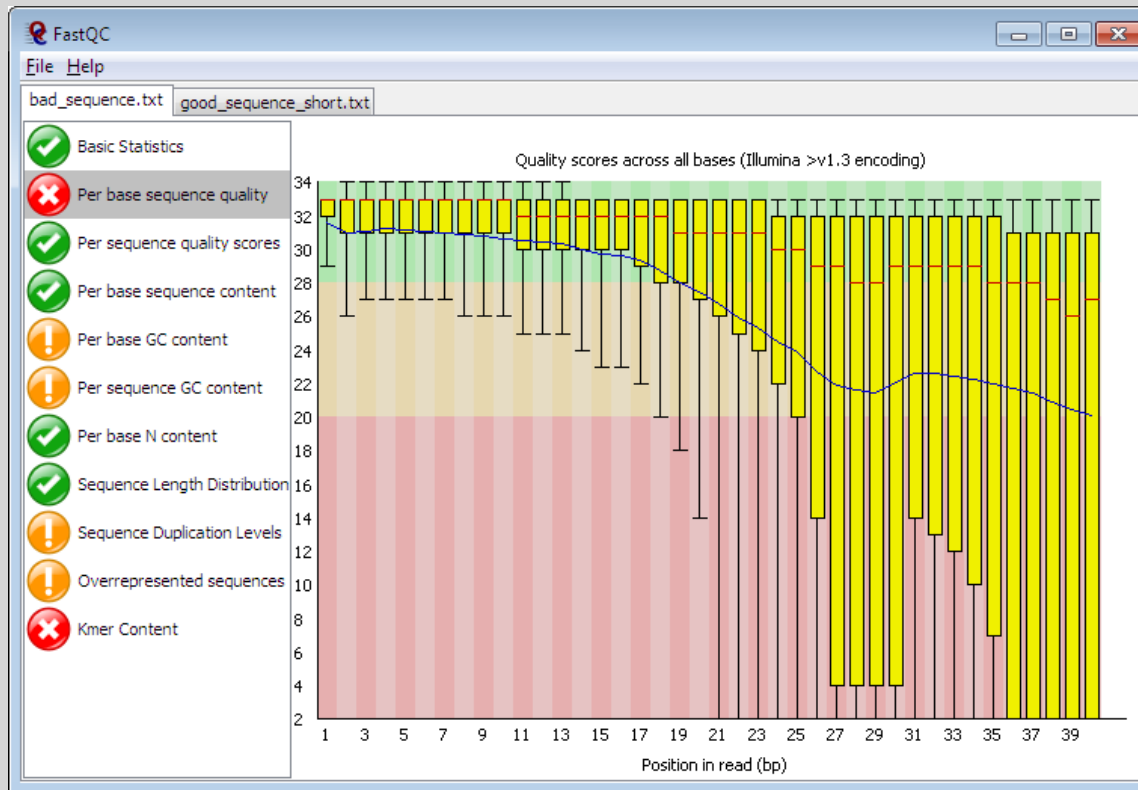




# FASTQC

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

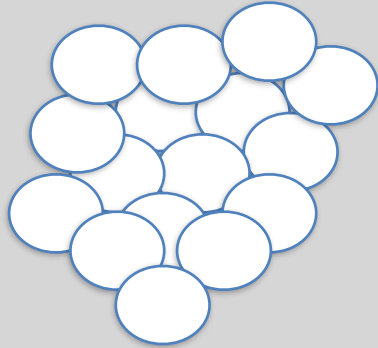
SSAHA

etc

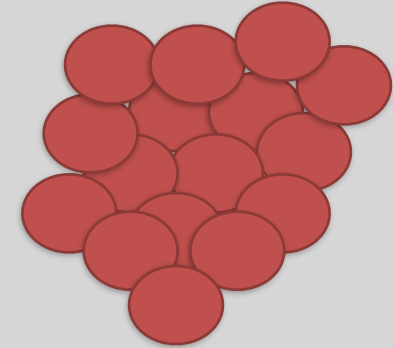
# RNA Sequencing

The absolute basics

Normal Cells

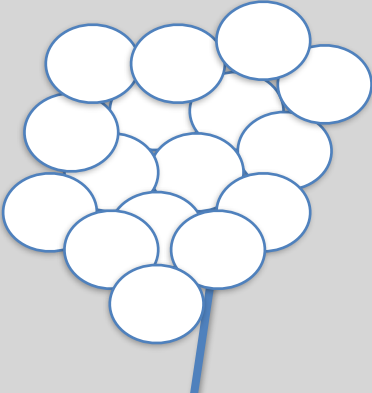


Mutated Cells

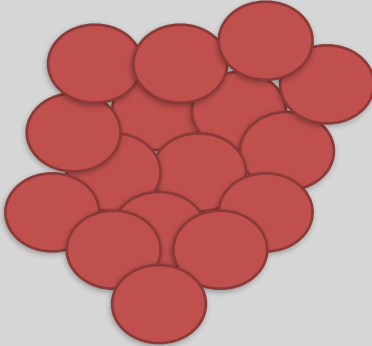


- The **mutated cells** behave differently than the **normal cells**
- We want to know what genetic mechanism is causing the difference
- One way to address this is to examine differences in gene expression via RNA sequencing...

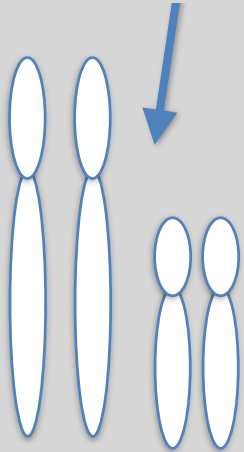
Normal Cells



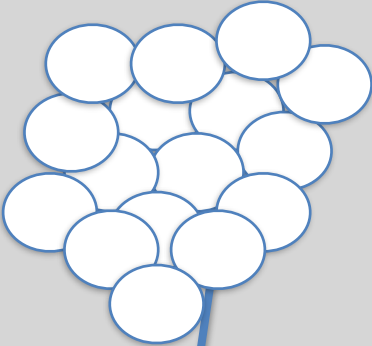
Mutated Cells



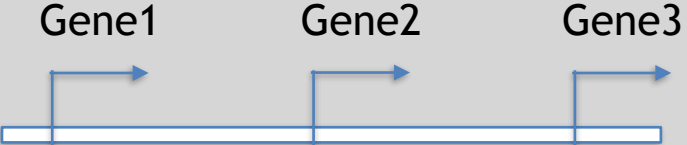
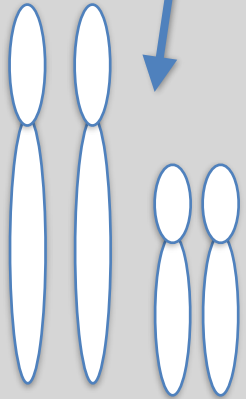
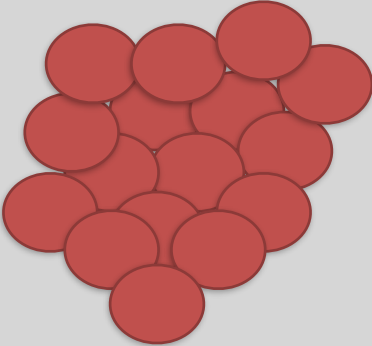
Each cell has a bunch of chromosomes



Normal Cells

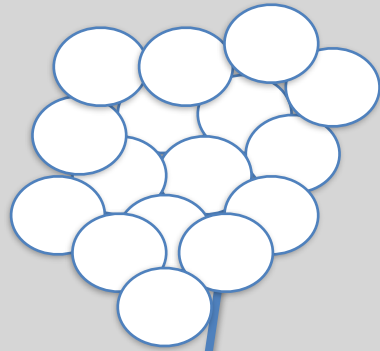


Mutated Cells

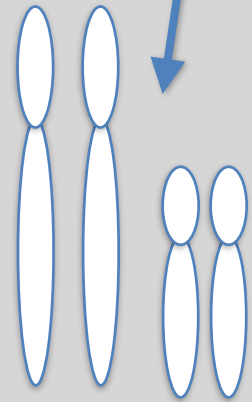
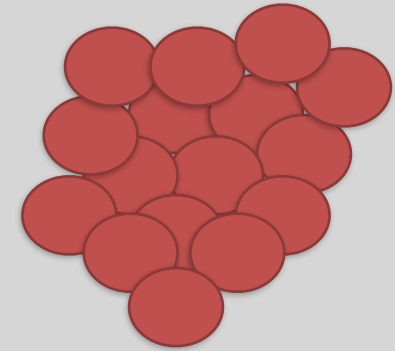


Each chromosome has a bunch of genes

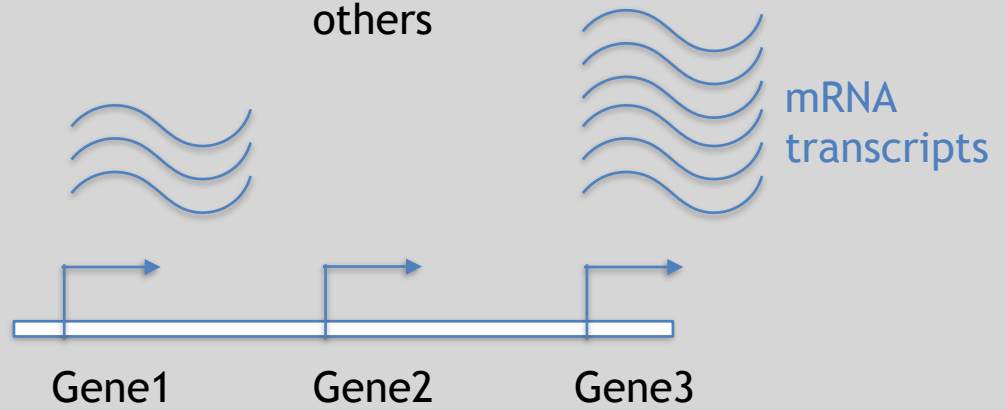
Normal Cells



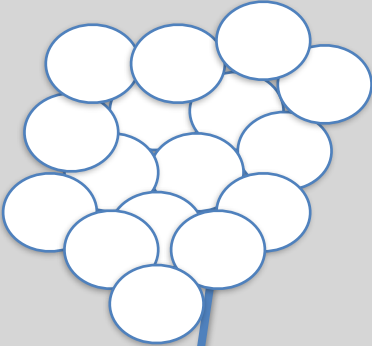
Mutated Cells



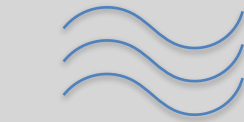
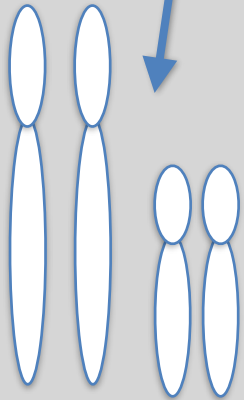
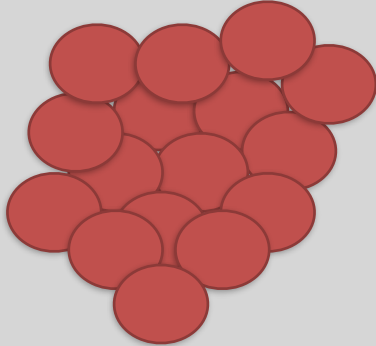
Some genes are active more than others



Normal Cells



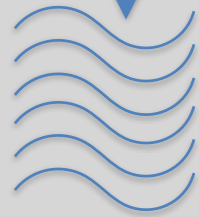
Mutated Cells



Gene 2 is not active



Gene 3 is the most active



mRNA transcripts



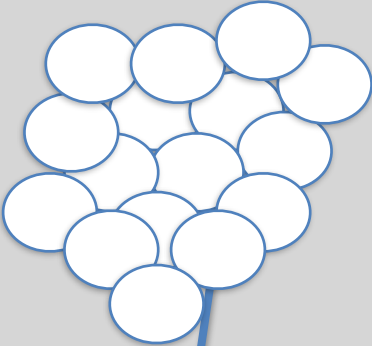
Gene1

Gene2

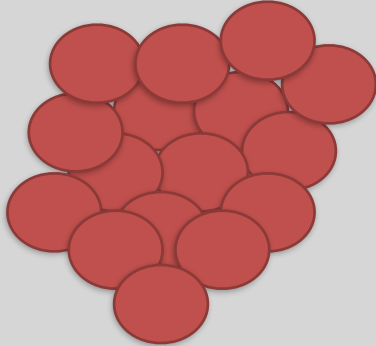
Gene3



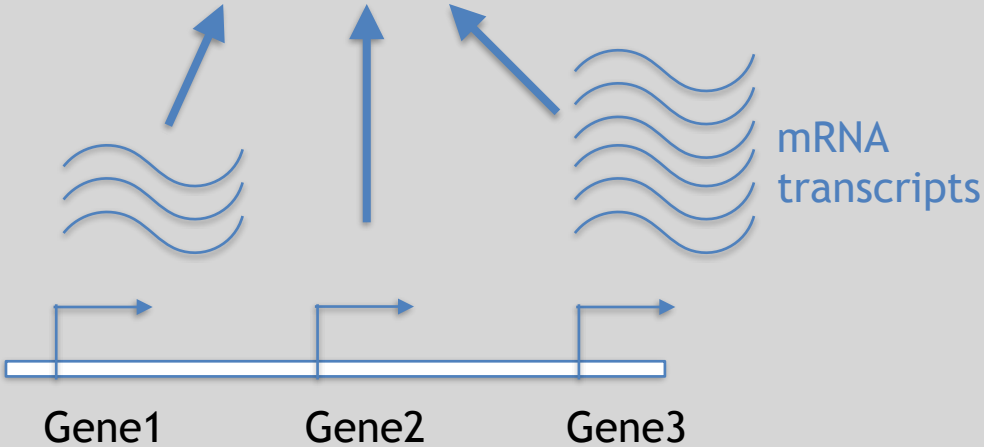
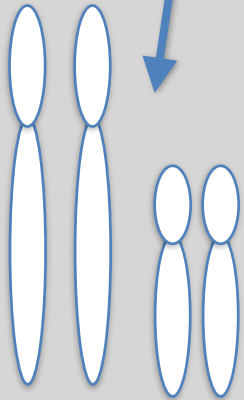
Normal Cells



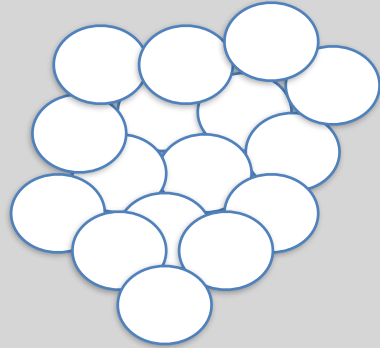
Mutated Cells



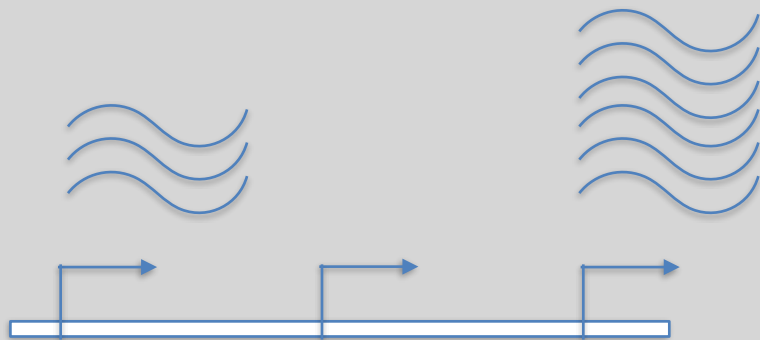
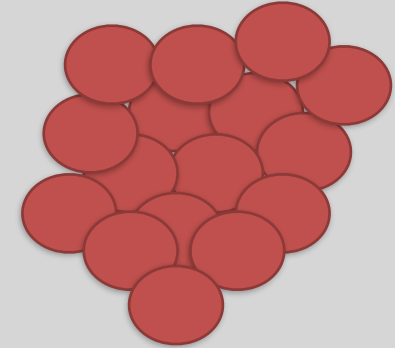
HTS tells us which genes are active, and how much they are transcribed!



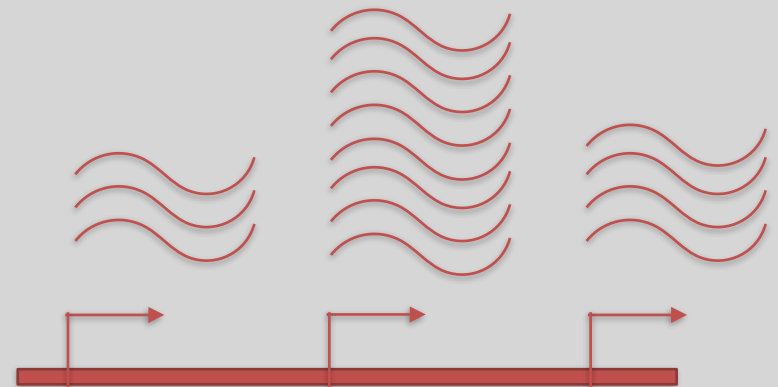
Normal Cells



Mutated Cells

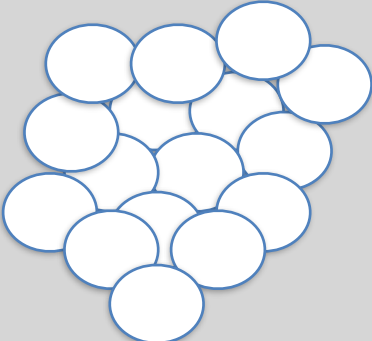


We use RNA-Seq to measure gene expression in normal cells ...

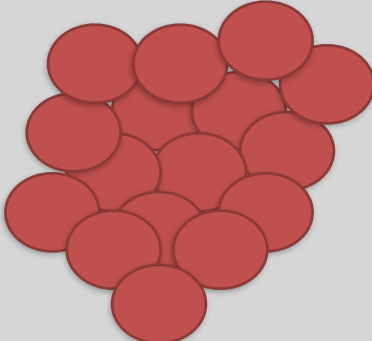


... then use it to measure gene expression in mutated cells

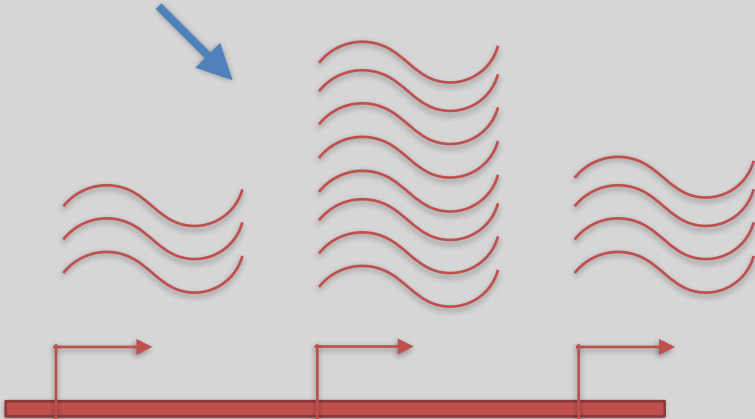
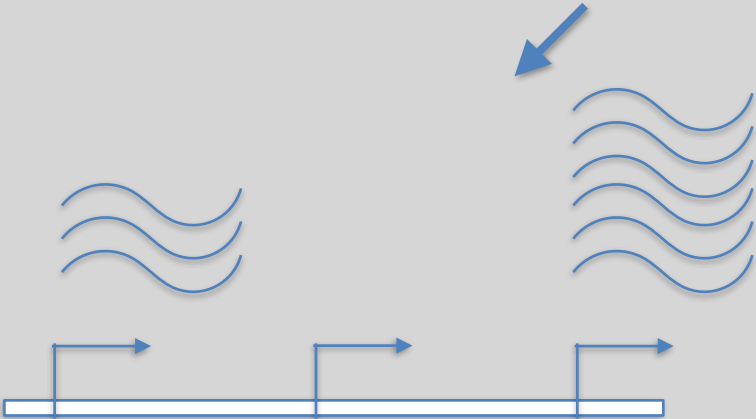
Normal Cells



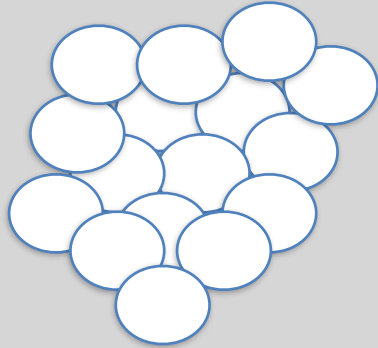
Mutated Cells



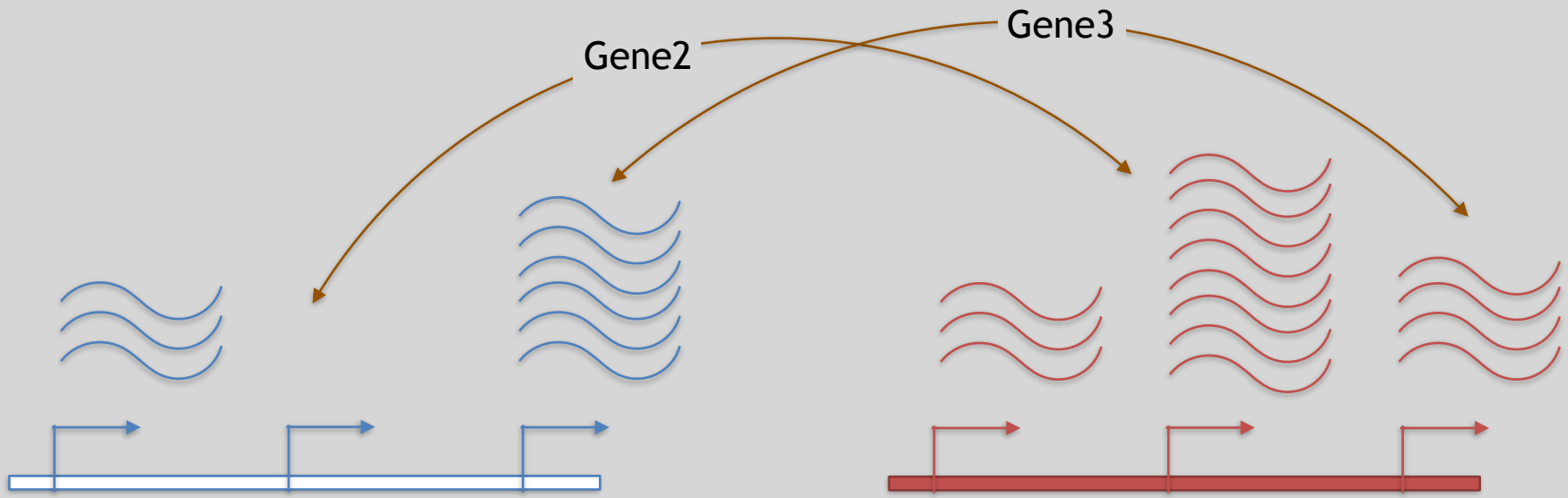
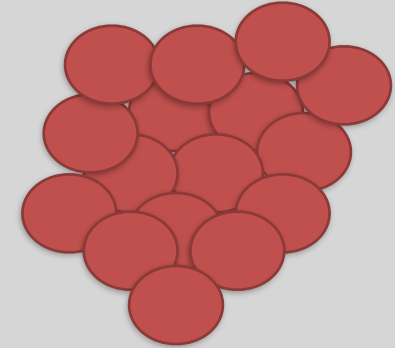
Then we can compare the two cell types to figure out what is different in the mutated cells!



Normal Cells



Mutated Cells



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

# 3 Main Steps for RNA-Seq:

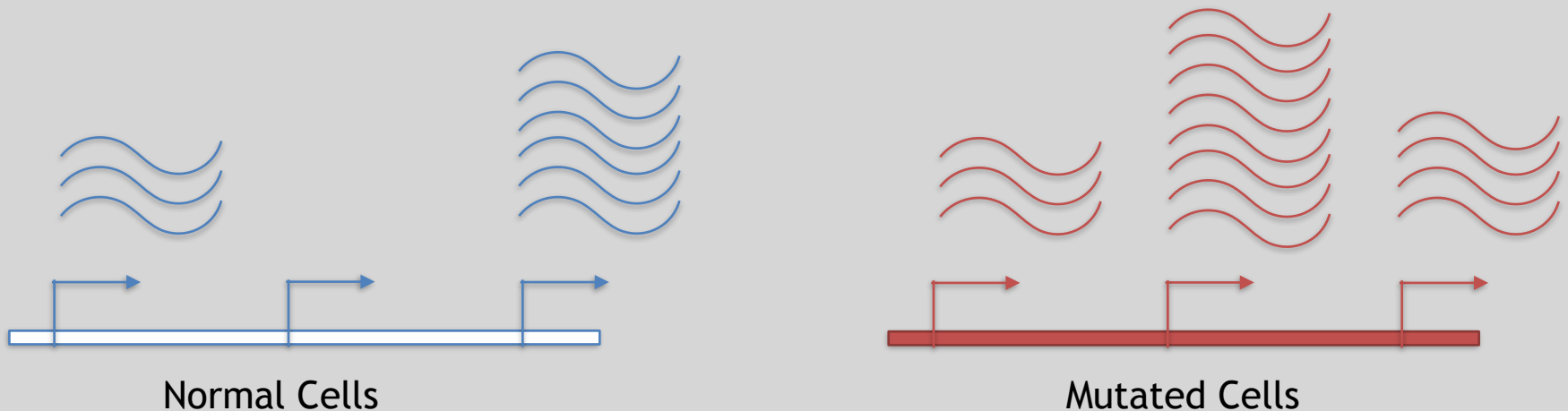
- 1) Prepare a sequencing library**  
(RNA to cDNA conversion via reverse transcription)
- 2) Sequence**  
(Using the same technologies as DNA sequencing)
- 3) Data analysis**  
(Often the major bottleneck to overall success!)

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/bgggn213\\_S19/lectures/#13](https://bioboot.github.io/bgggn213_S19/lectures/#13)

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

Feedback:

[Muddy Point Assessment]



# **Additional Reference Slides** on SAM/BAM Format and Sequencing Methods

# Sequence Alignment

Reference

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

# SAM Utilities

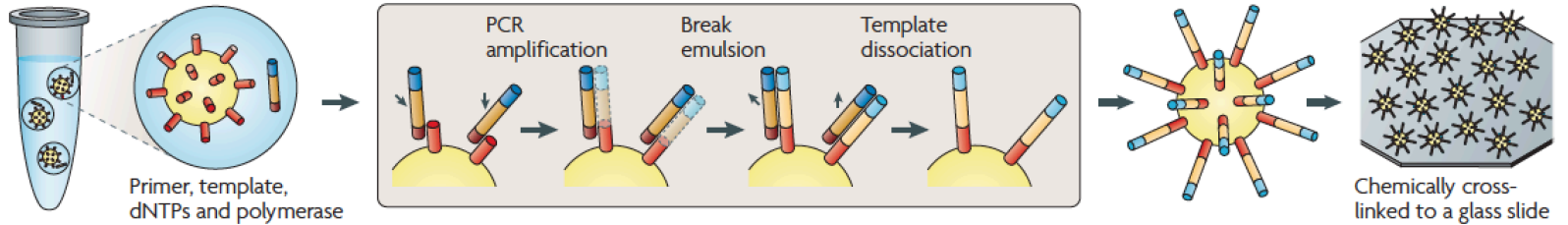
- **Samtools** is a common toolkit for analyzing and manipulating files in SAM/BAM format
  - <http://samtools.sourceforge.net/>
- **Picard** is a another set of utilities that can used to manipulate and modify SAM files
  - <http://picard.sourceforge.net/>
- These can be used for viewing, parsing, sorting,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

# Additional Reference Slides on Sequencing Methods

# Roche 454 - Pyrosequencing

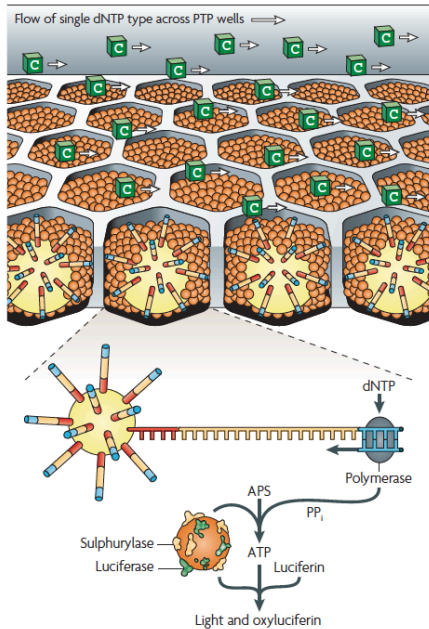
## a Roche/454, Life/APG, Polonator Emulsion PCR

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



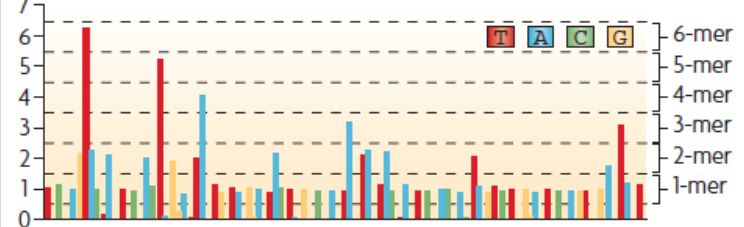
## c Roche/454 — Pyrosequencing

1-2 million template beads loaded into PTP wells



## d Flowgram

TCAGGTTTTTTTAAACAATCAACTTTTTGGATTAAAAATGTAGATAACTG  
CATAAATTAATAACATCACATTAGTCTGATCAGTGAATTTAT

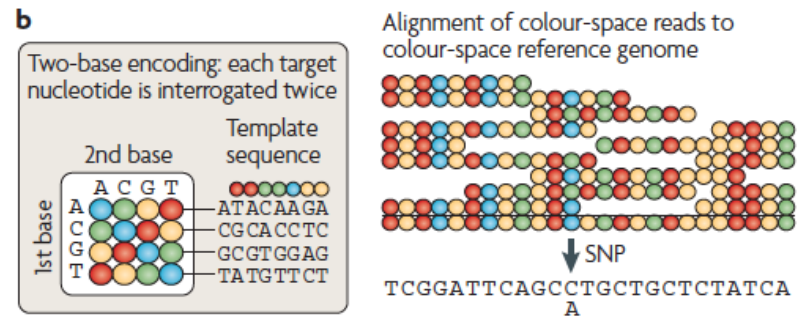
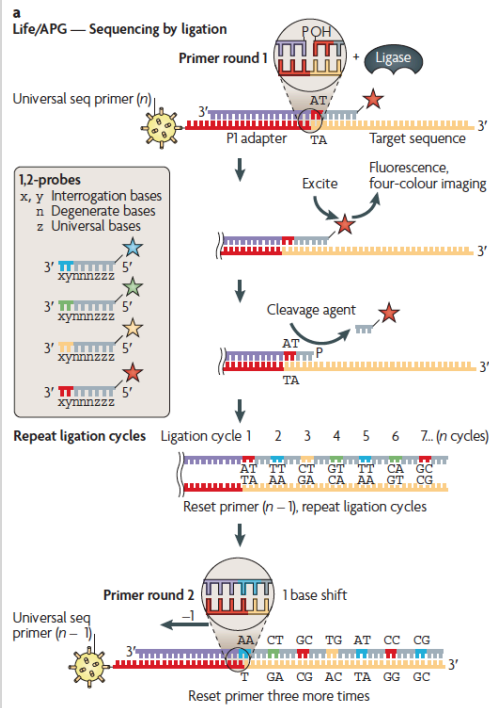
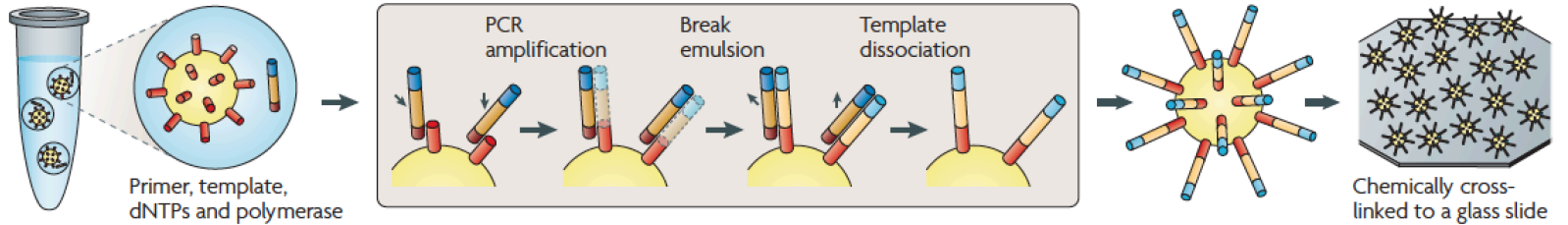




# Life Technologies SOLiD - Sequence by Ligation

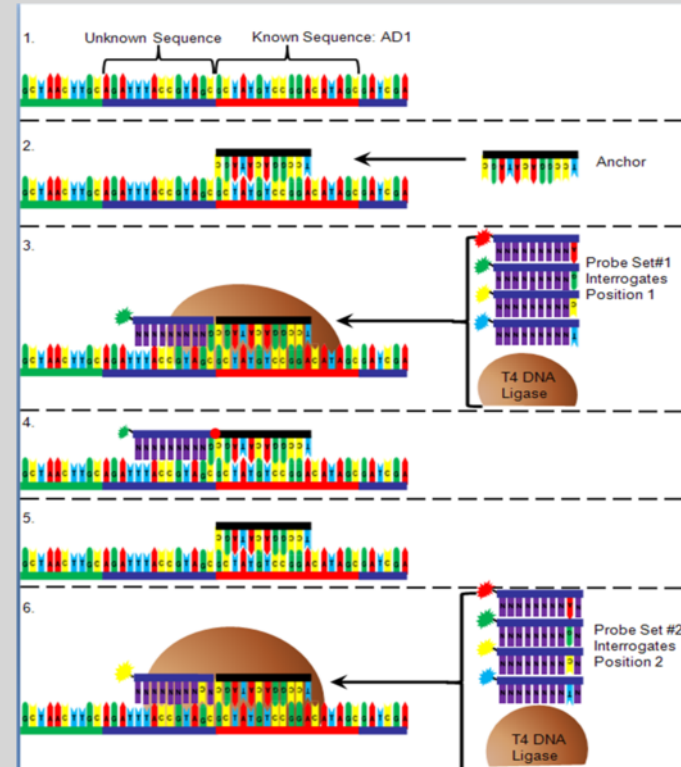
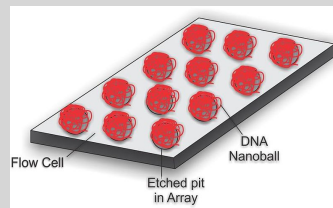
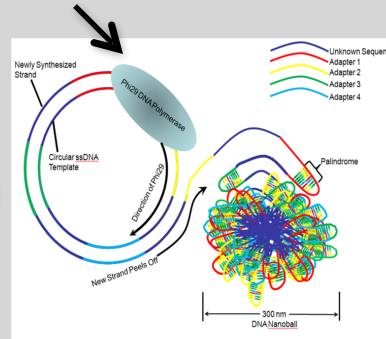
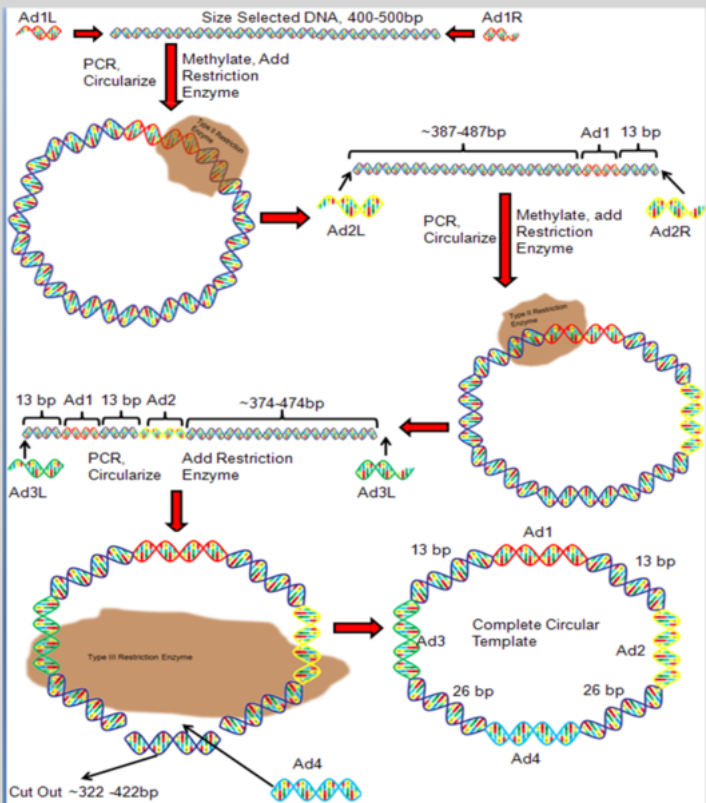
## a Roche/454, Life/APG, Polonator Emulsion PCR

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



# Complete Genomics - Nanoball Sequencing

Has proofreading ability!

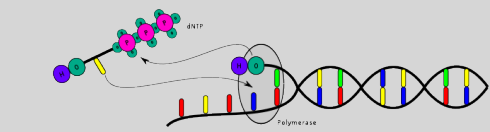


# “Benchtop” Sequencers

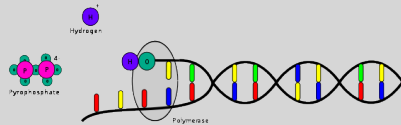
- Lower cost, lower throughput alternative for smaller scale projects
- Currently three significant platforms
  - Roche 454 GS Junior
  - Life Technology Ion Torrent
    - Personal Genome Machine (PGM)
    - Proton
  - Illumina MiSeq

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 <sup>a,b</sup>	\$225 <sup>c</sup>	10 Mb (100 bases)	3 h	\$22.5	3.3
(316 chip)		\$425	100 Mb <sup>d</sup> (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

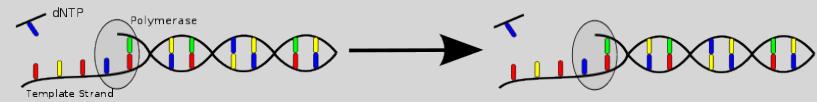
# PGM - Ion Semiconductor Sequencing



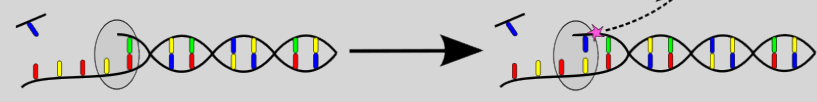
Polymerase integrates a nucleotide.



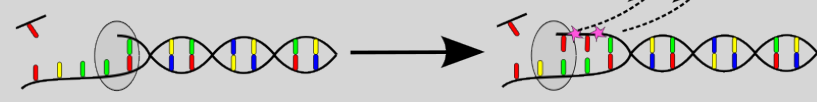
Hydrogen and pyrophosphate are released.



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.

