

BIMM 143
Genome Informatics I
Lecture 13
Barry Grant
UC San Diego
<http://thegrantlab.org/bimm143>

TODAYS MENU:

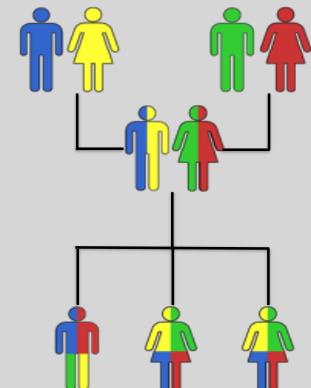
- **What is a Genome?**
 - Genome sequencing and the Human genome project
- **What can we do with a Genome?**
 - Compare, model, mine and edit
- **Modern Genome Sequencing**
 - 1st, 2nd and 3rd generation sequencing
- **Workflow for NGS**
 - RNA-Sequencing and Discovering variation

Genetics and Genomics

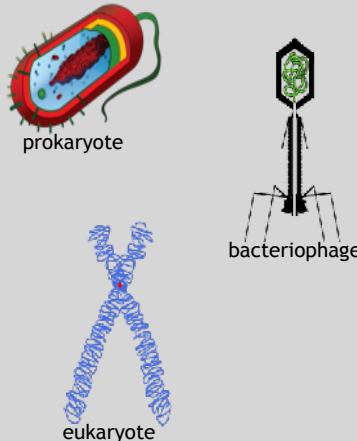
- **Genetics** is primarily the study of individual genes, mutations within those genes, and their inheritance patterns in order to understand specific traits.
- **Genomics** expands upon classical genetics and considers aspects of the entire genome, typically using computer aided approaches.

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



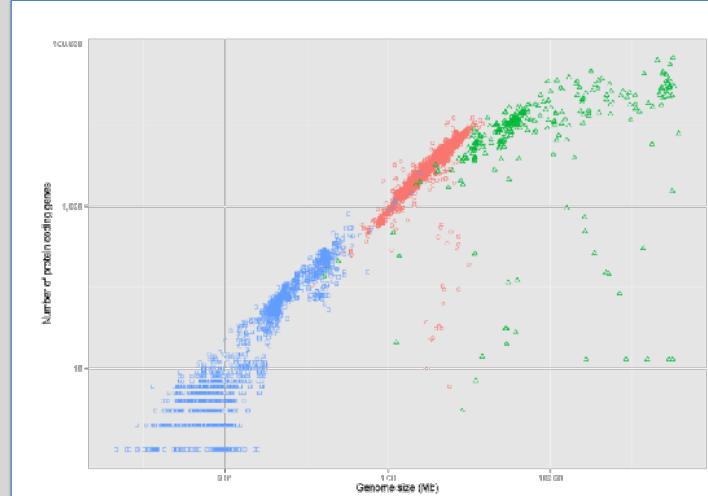
Genomes come in many shapes



- 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)

Prokaryote by [Mariela Ruiz Villareal](#) | Bacteriophage image by [Salome](#) / CC BY-SA | Eukaryote image by [Manuel Moncke](#) / CC BY-SA

Genomes come in many sizes

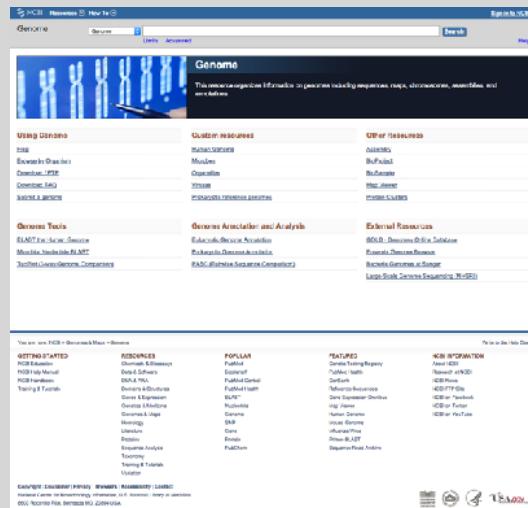


Modified from image by [Eduardo](#) / CC BY-SA

Genome Databases

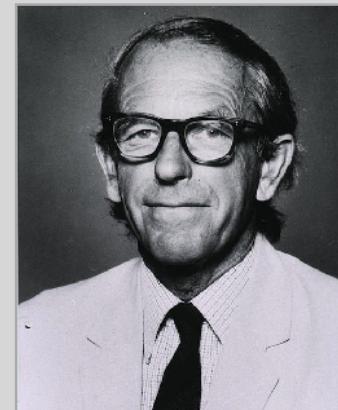
NCBI Genome:

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



This screenshot shows the homepage of the NCBI Genome database. The top navigation bar includes links for "Home", "Search", "Help", and "Log In". The main content area features a large image of a karyogram and the text: "This resource organizes information on genomes including sequences, maps, chromosomes, assemblies, and gene delivery". Below this are several sections: "Using Genome", "Genome Tools", "Genome Databases", "Genome Resources", "Genome Analysis and Analysis", and "External Resources". At the bottom, there are links for "Help", "Feedback", "Log In", and "Logout".

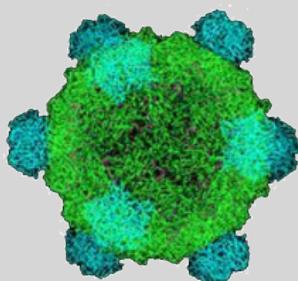
Early Genome Sequencing



http://en.wikipedia.org/wiki/Frederick_Sanger

- 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 ~1 / 10000 bases

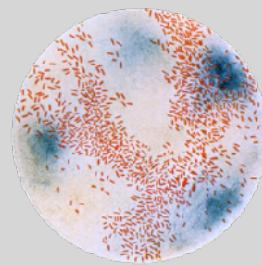
The First Sequenced Genomes



Bacteriophage φ-X174

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

http://en.wikipedia.org/wiki/Phi_X_174



Haemophilus influenzae

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

<http://phl.cdc.gov/>

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



Jane Ades, Courtesy: National Human Genome Research Institute

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

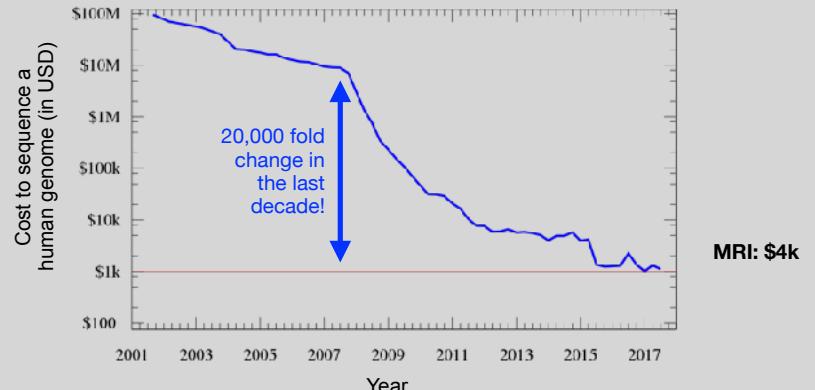
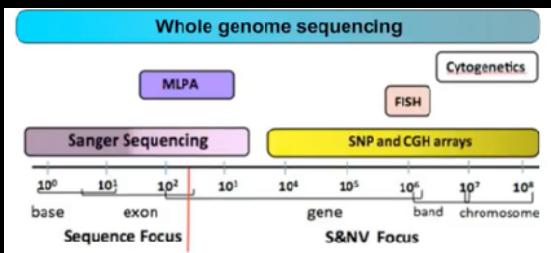


Image source: https://en.wikipedia.org/wiki/Carlson_curve

Whole genome sequencing transforms genetic testing



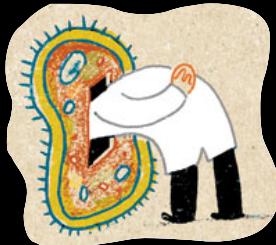
- 1000s of single gene tests
- Structural and copy number variation tests
- Permits hypothesis free diagnosis

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
Protein expression	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

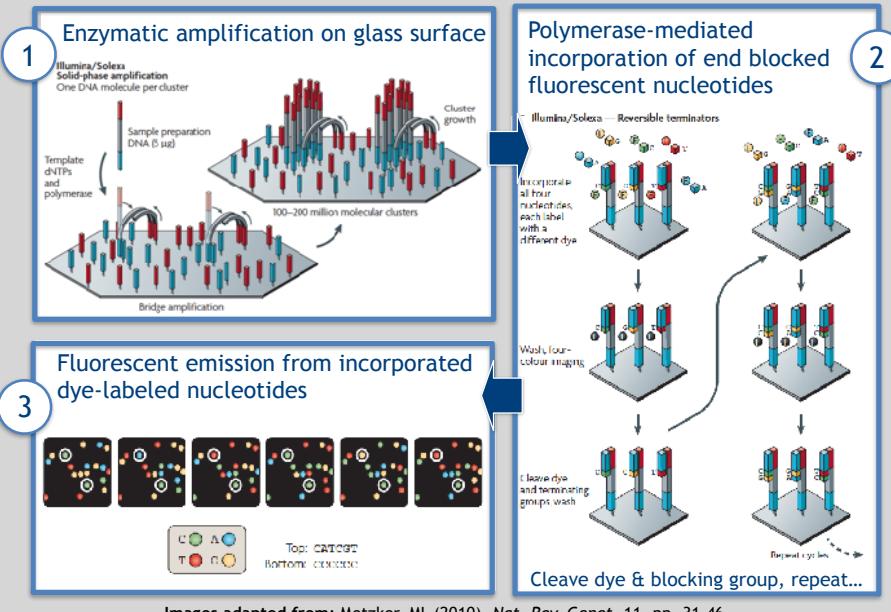
- **Sequencing by Synthesis:** Uses a polymerase to incorporate and assess nucleotides to a primer sequence
 - 1 nucleotide at a time
- **Sequencing by Ligation:** Uses a ligase to attach hybridized sequences to a primer sequence
 - 1 or more nucleotides at a time (e.g. dibase)

Modern NGS Sequencing Platforms

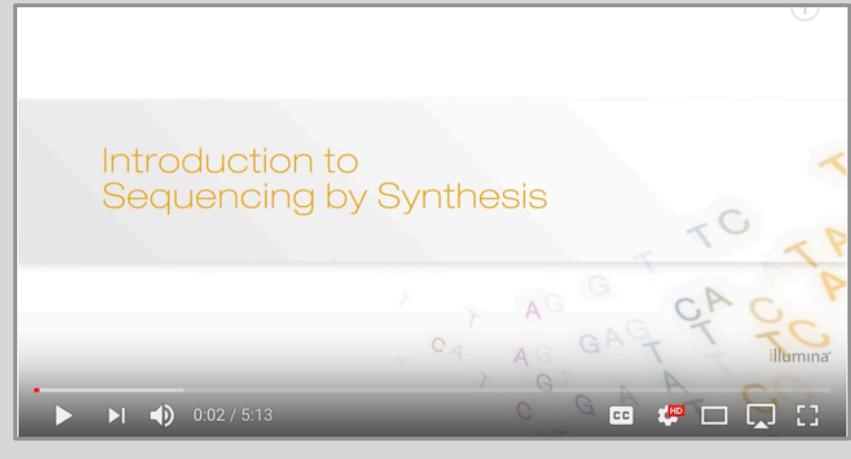
	Roche/454	Life Technologies SOLiD	Illumina Hi-Seq 2000
Library amplification method	emPCR® on bead surface	emPCR® on bead surface	Enzymatic amplification on glass surface
Sequencing method	Polymerase mediated incorporation of unlabelled nucleotides	Ligase mediated addition of 2-base encoded fluorescent 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

Modified from Mardis, ER (2011), Nature, 470, pp. 198-203

Illumina - Reversible terminators

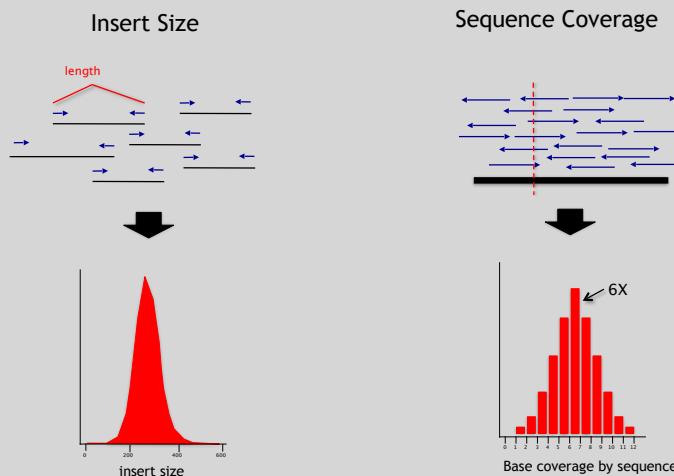


Illumina Sequencing - Video



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

NGS Sequencing Terminology



Summary: “Generations” of DNA Sequencing

	First generation	Second generation ^b	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 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 base	Low-to-moderate cost per base
	Low cost per run	High cost per run	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

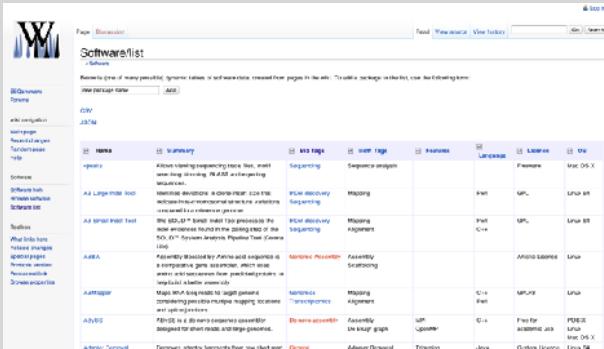
Schadt, EE et al (2010), *Hum. Mol. Biol.*, 19(R12), pp. R227-R240

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!

SeqAnswers Wiki

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



Tool	Summary	Tags	Type	Version	Features	Licenses	License	OS
agpeptid	Automatically reporting peptide lists, even after long sequencing. It also distinguishes between peptides.	Peptide discovery	Sequence analysis	0.1.0	Peptide mapping	GPL	GPL	Mac OS X
All Large PSM tool	Identifying peptides. It finds many PSMs that are not found in other reference databases.	Peptide discovery	Mapping	0.1.0	Peptide sequencing	GPL	GPL	Linux/OSX
All Small PSM tool	The small PSM tool is designed for the small projects that don't need the NCE, DTA, or Peptide Tree Coverage tool.	Peptide discovery	Peptide alignment	0.1.0	Peptide sequencing	GPL	GPL	Linux/OSX
ABBA	Perfomrantly classified by ABBA and separated in a comparative gene assembler, which is a new way to assemble genomes. This tool is able to handle assembly of large datasets.	NanoReads	Assembly	0.1.0	Assembly	Apache License	Apache	
Autospacer	Map MAQ output to target genome and compare it to reference genome, then map spacers to target genome.	NanoReads	Motif alignment	0.1.0	MAQ	GPL	GPL	Linux
AVES	PDFC is a bioinformatics assembler designed for short reads and large genomes.	De novo assembly	Frequently De Novo	0.1.0	De Novo assembly	PBS	PBS	Mac OS X

The first direct RNA sequencing by nanopore

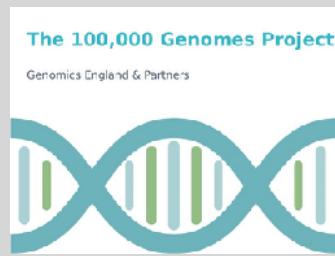
- For example this new nanopore sequencing method was just published!
<https://www.nature.com/articles/nmeth.4577>
- "Sequencing the RNA in a biological sample can unlock a wealth of information, including the identity of bacteria and viruses, the nuances of alternative splicing or the transcriptional 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."

Side-Note:

What can we do with all this sequence information?

Population Scale Analysis

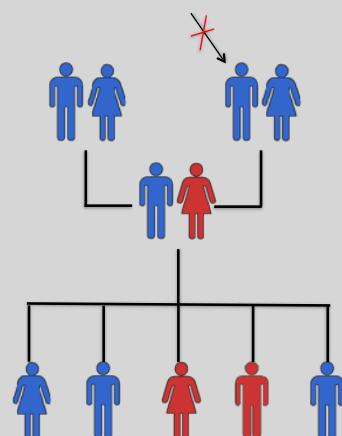
We can now begin to assess genetic differences on a very large scale, both as naturally occurring variation in human and non-human populations as well somatically within tumors



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

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



“Variety’s the very spice of life”

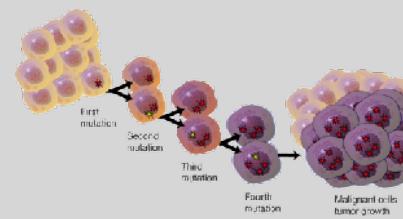
-William Cowper, 1785

“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.

Somatic Variation

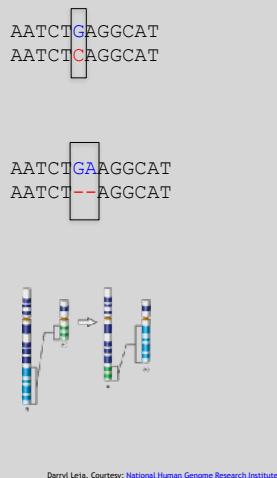


- Mutations in non-germline cells that are not passed along to offspring
- Can occur during mitosis or from the environment itself
- Are an integral part in tumor progression and evolution

Daryl Leja, Courtesy: National Human Genome Research Institute

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



Discovering Variation: SNPs and INDELs

SNP
sequencing error or genetic variant?
sequencing error or genetic variant?
INDEL

reference genome

ATCCTGATTCCGTGAACGTTATCGACCATCCGATCGA	TTATCGACATCCGATCGAACTGTCA CGGGCAAGCTGATCG
ATCCTGATTCCGTGAACGTTATCGACCATCCGATCGA	TCGACGATCCGATCGAACTGTCA CGGGCAAGCTGATCG
CCGTGAACGTTATCGACCATCCGATCGAACTGTCA CGGC	ATCCGATCGAACTGTCA CGGGCAAGCTGATCG CGAT
GGTGAACTGTTATCGACCATCCGATCGAACTGTCA CGCG	TCCGACGTTATCGACCATCCGATCGAACTGTCA CGATCG
TGAACGTTATCGACCATCCGATCGAACTGTCA CGGC	TCCGATCGAACTGTCA CGGGCAAGCTGATCG CGATCG
GTTATCGACCATCCGATCGAACTGTCA CGGGCAAGCT	TGTCAGGGCAAGCTGATCG CGATCGATGCTAGTG
TTATCGACCATCCGATCGAACTGTCA CGGGCAAGCT	TCAGCGGCAAGCTGATCG CGATCGATGCTAGTG

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

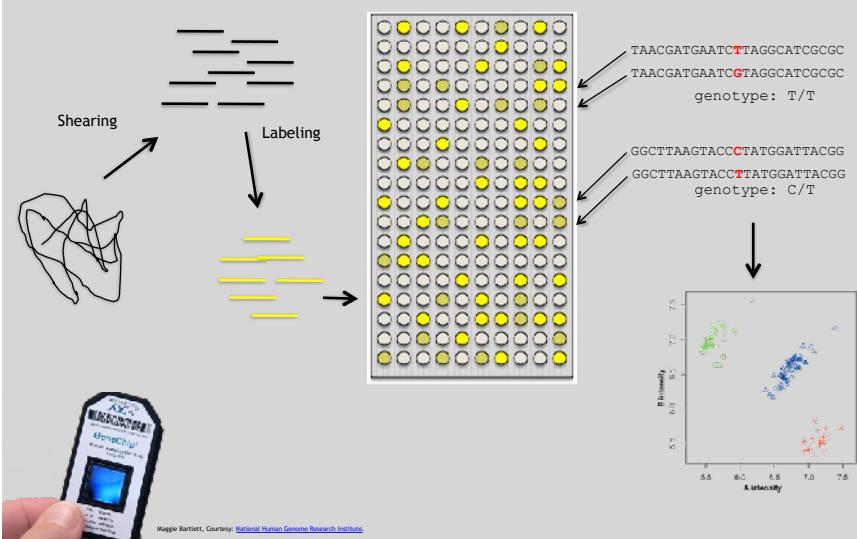
Numbers change depending on ancestry!

[Numbers from: 1000 Genomes Project, Nature, 2012]

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

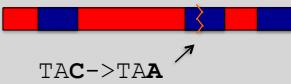
SNP Microarrays



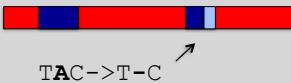
Impact of Genetic Variation

There are numerous ways genetic variation can exhibit functional effects

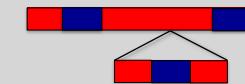
Premature stop codons



Frameshift mutation



Gene or exon deletion



Transcription factor binding disruption

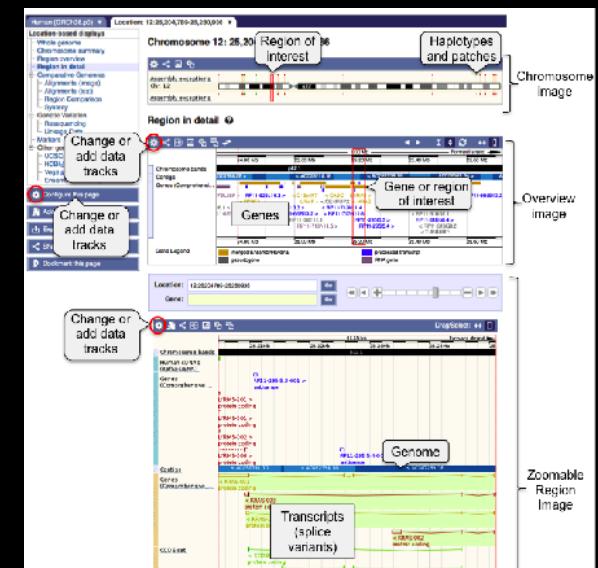


Hand-on time!

https://bioboot.github.io/bimm143_S18/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>



Access a jetstream galaxy instance!

Use assigned IP address

Do it Yourself!

The screenshot shows the Galaxy web interface with a workflow titled "Bowtie2 - map reads against reference genome (Galaxy Version 2.2.6.2)". The workflow consists of several steps:

- Step 1: "Is this single or paired library" (Single-end)
- Step 2: "FASTQ file" (Must be of datatype "Fastq")
- Step 3: "Write unaligned reads (in fastq format) to separate file(s)" (Yes)
- Step 4: "Will you select a reference genome from your history or use a built-in index?" (Use a built-in genome index)
- Step 5: "Select reference genome" (Buildin (hg19))
- Step 6: "Set read groups information?" (Do not set)
- Step 7: "Specifying read group information can greatly simplify your downstream analyses by allowing combining multiple datasets."
- Step 8: "Select analysis mode" (Default setting only)
- Step 9: "Do you want to use presets?" (No, just use defaults)
- Step 10: "Very fast end-to-end (-very-fast)" (radio button selected)
- Step 11: "Fast end-to-end (-fast)"
- Step 12: "Sensitive end-to-end (-sensitive)"
- Step 13: "Very sensitive end-to-end (-very-sensitive)"
- Step 14: "Very fast local (-very-fast-local)"
- Step 15: "Fast local (-fast-local)"
- Step 16: "Sensitive local (-sensitive-local)"
- Step 17: "Very sensitive local (-very-sensitive-local)"

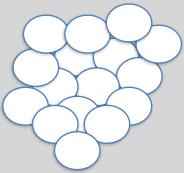
The history panel on the right shows several completed steps, including:

- 23: "Index on data 18 and data 17" (Completed)
- 24: "Index on data 18 and data 17" (Completed)
- 25: "Cufflinks on data 18 and data 16. Skipped transcripts" (Completed)
- 26: "Cufflinks on data 18 and data 16. assembled transcripts" (Completed)
- 27: "Cufflinks on data 18 and data 16. transcript expression" (Completed)
- 28: "Cufflinks on data 18 and data 16. gene expression" (Completed)
- 29: "Index on database hg19" (Completed)
- 30: "cufflinks v3.2.1 cufflinks -q --no-update-check -I 300000 -F 0.100000 -J 0.100000 -p 2 -z /tmp/cufflinks_index_4 /opt/galaxy/galaxy-apps/database/datasets/0000/dataset_4" (Completed)

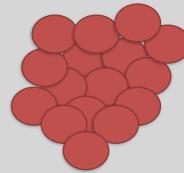
RNA Sequencing

The absolute basics

Normal Cells

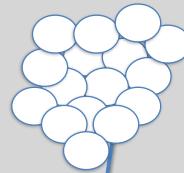


Mutated Cells

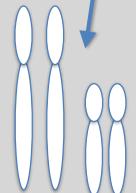


- The **mutated cells** behave differently than the **normal cells**
- We want to know what genetic mechanism is causing the difference
- One way to address this is to examine differences in gene expression via RNA sequencing...

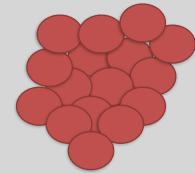
Normal Cells

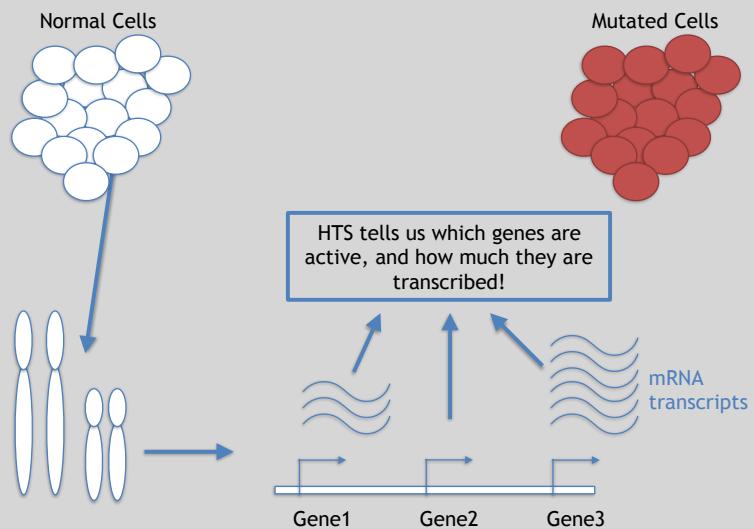
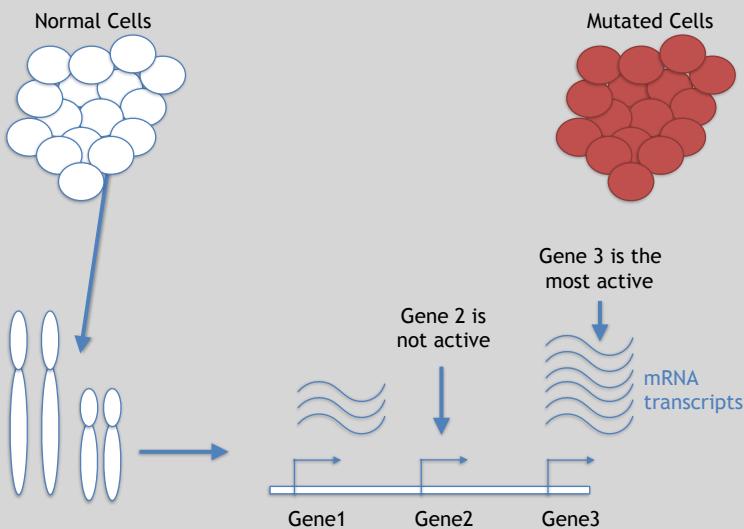
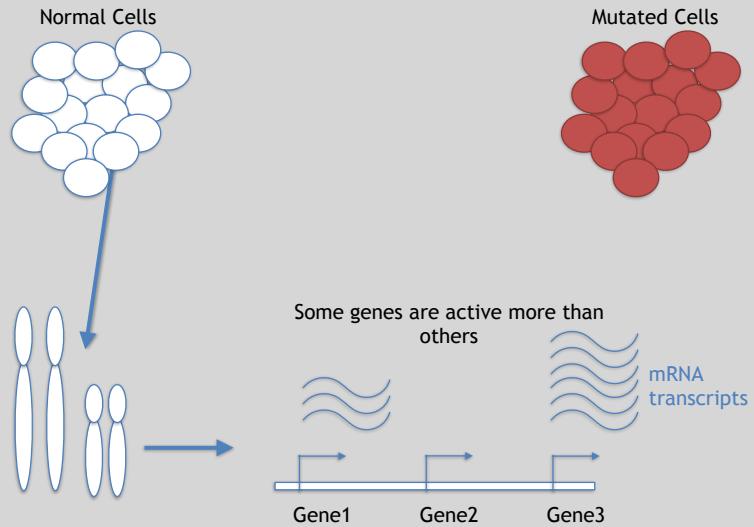
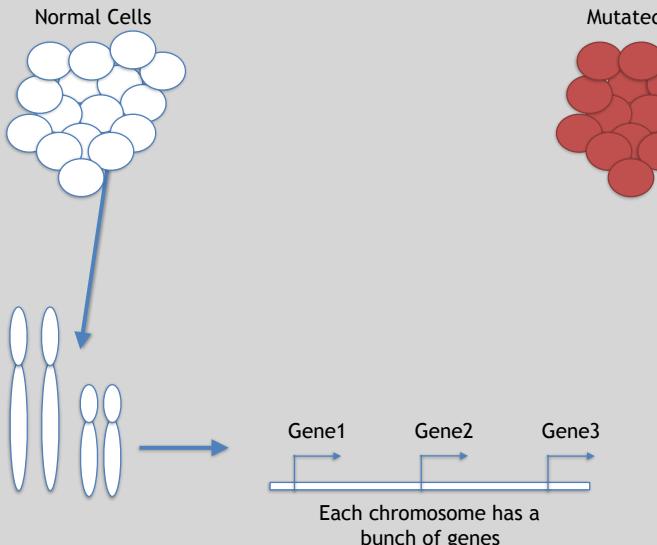


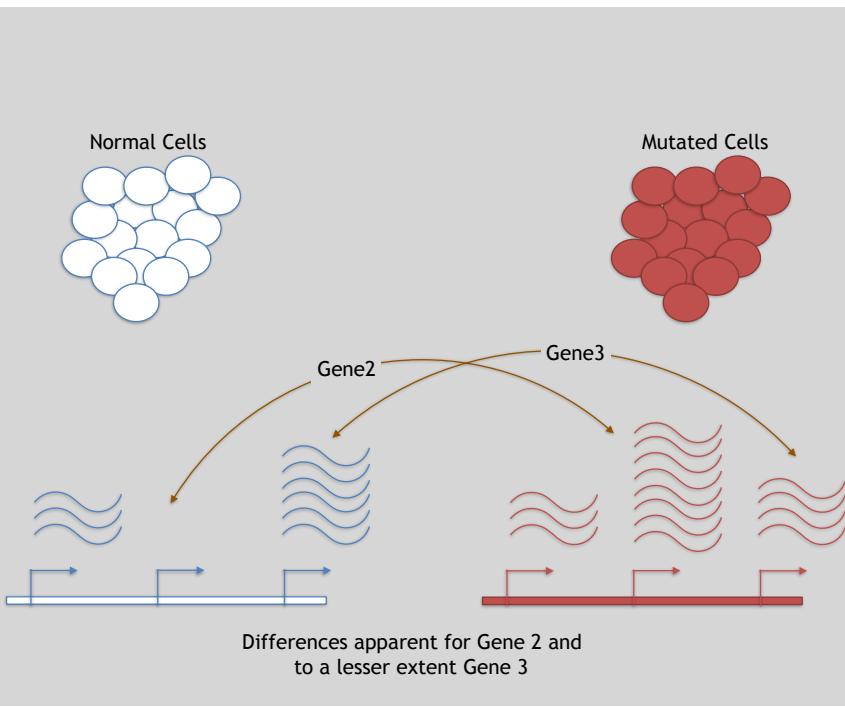
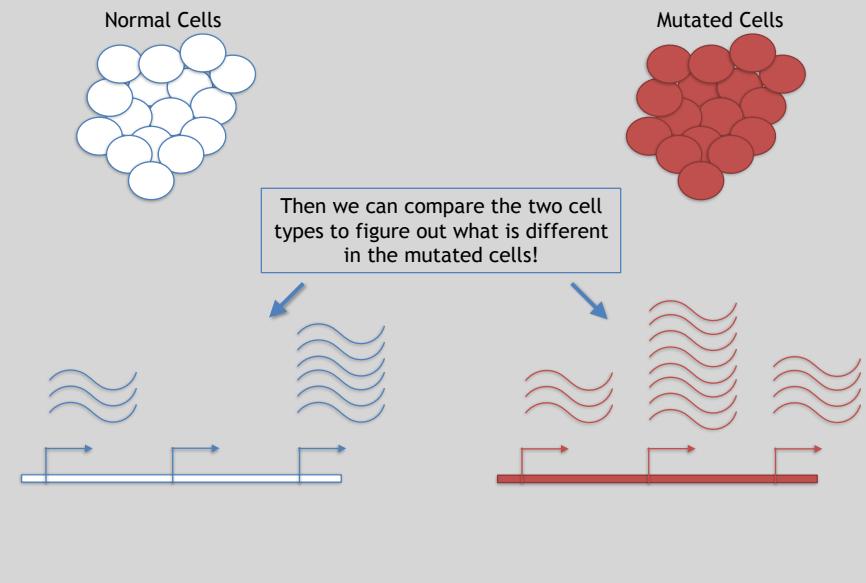
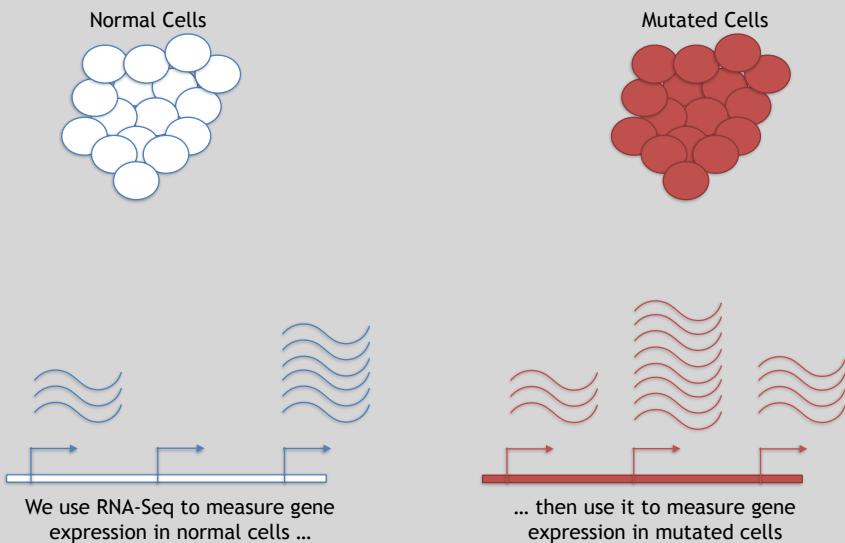
Each cell has a bunch of chromosomes



Mutated Cells







3 Main Steps for RNA-Seq:

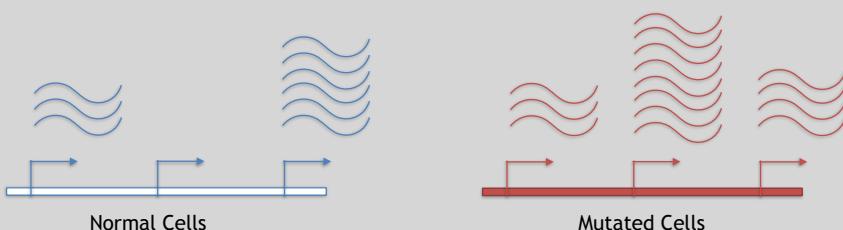
- 1) Prepare a sequencing library**
(RNA to cDNA conversion via reverse transcription)
- 2) Sequence**
(Using the same technologies as DNA sequencing)
- 3) Data analysis**
(Often the major bottleneck to overall success!)

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



Additional Reference Slides

(On FASTQ format, ASCII Encoded Base Qualities, FastQC, Alignment and SAM/BAM formats)

Hands-on worksheet:
https://bioboot.github.io/bimm143_W18/lectures/#13

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

Raw data usually in FASTQ format

```
@NS500177:196:HFTTAFXX:1:11101:10916:1458 2:N:0:CGGGCTG  
ACACGACATGAGGTGACAGTCACGGAGATAAGATCAATGCCCTCATTAAGCAGCCGGTAA  
+  
AAAAAEEEEEEEEE//AAAAAEEEEEEEEE//EE//<<EE/AAAFAEE//EEEAEAAE<  
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 forth line contains the quality scores for each base in the sequenced fragment (these are ASCII encoded...)

ASCII Encoded Base Qualities

```
@NS500177:196:HFTTAFXX:1:11101:10916:1458 2:N:0:CGGGCTG  
ACACGACGATGAGGTGACAGTCACGGAGGATAAGATCAATGCCCTCATTAAGCAGCCGGTAA  
+  
AAAAAEEEEEEEEE//AEEEAEeeeeeee/EE/<<EE/AEEEAEE//EEEAEAAE< 4
```

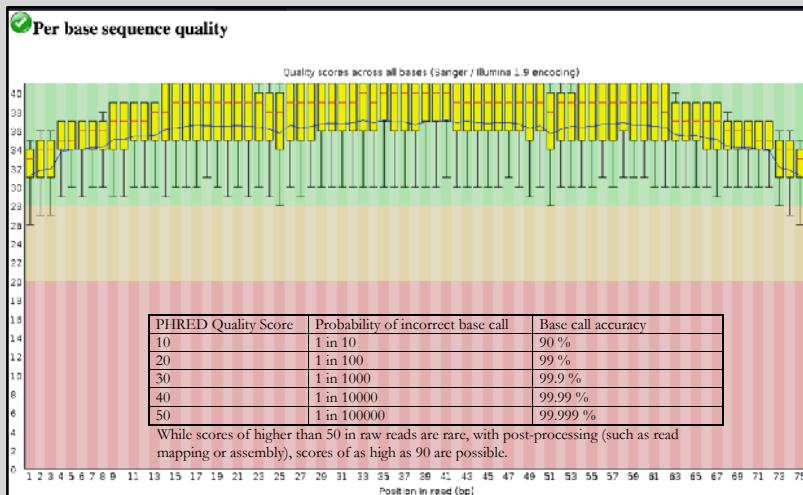
- Each sequence base has a corresponding numeric quality score encoded by a single ASCII character typically on the 4th line (see ④ above)
- ASCII characters represent integers between 0 and 127
- Printable ASCII characters range from 33 to 126
- Unfortunately there are 3 quality score formats that you 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("DDDCDEDCCDDDBBDDCC@") ) - 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 35 33 33 35 35 35 34 34 31
> prob <- 10**(-phred/10)
```

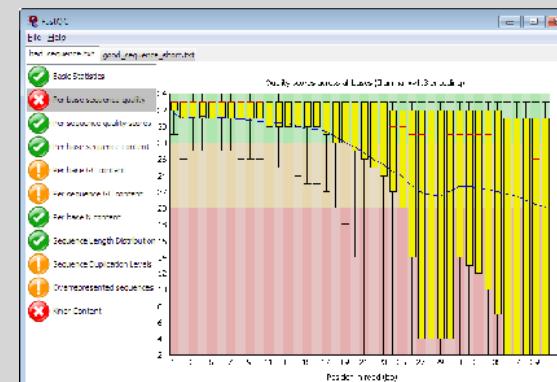
FastQC Report



FASTQC

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

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



SAM Utilities

- Samtools is a common toolkit for analyzing and manipulating files in SAM/BAM format
 - <http://samtools.sourceforge.net/>
- Picard is 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)

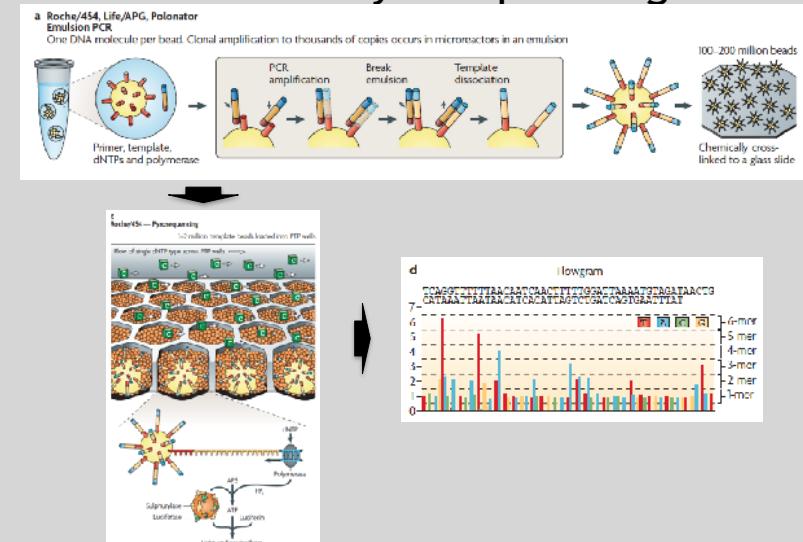
Additional Reference Slides
on Sequencing Methods

Do it Yourself!

Genome Analysis Toolkit (GATK)

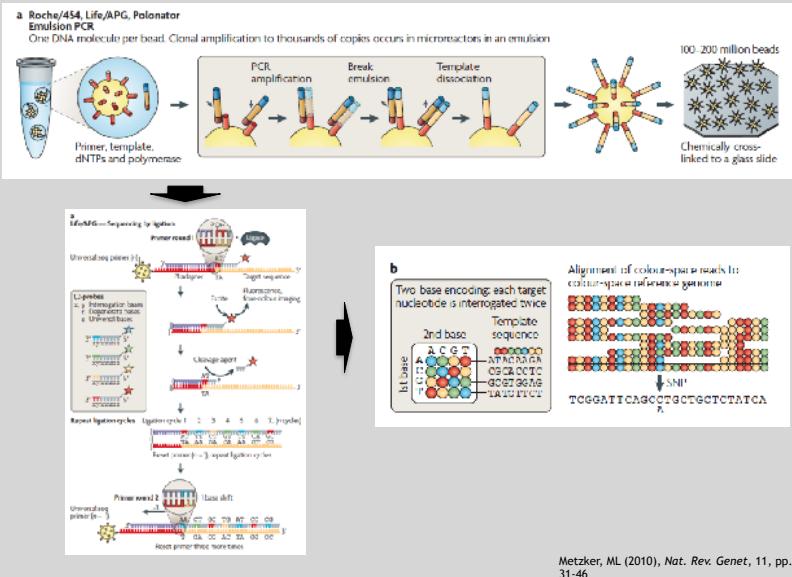
- Developed in part to aid in the analysis of 1000 Genomes Project data
- Includes many tools for manipulating, filtering, and utilizing next generation sequence data
- <http://www.broadinstitute.org/gatk/>

Roche 454 - Pyrosequencing

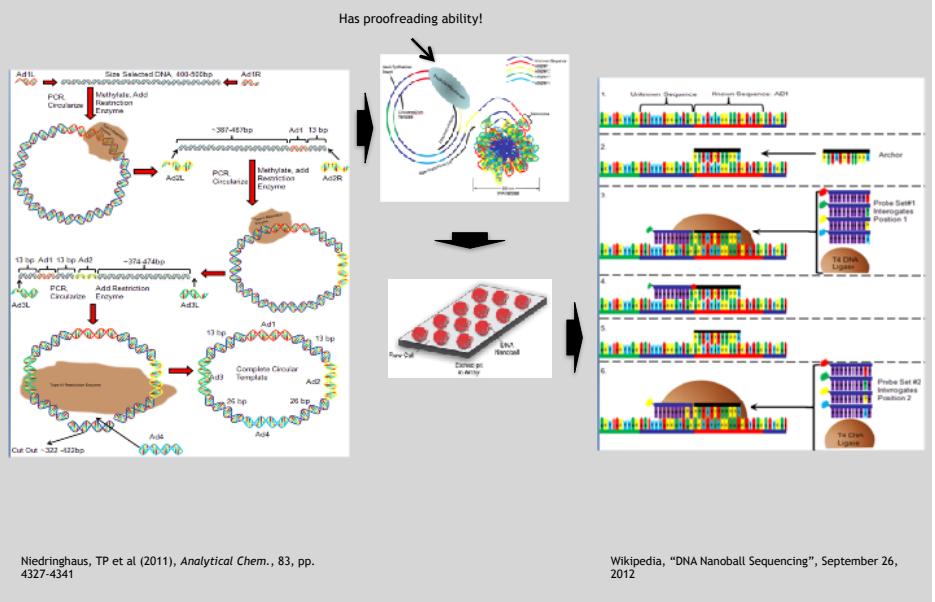


Metzker, ML (2010). *Nat. Rev. Genet.*, 11, pp. 31-46.

Life Technologies SOLiD - Sequence by Ligation



Complete Genomics - Nanoball Sequencing



"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

Loman, NJ (2012), *Nat. Biotech.*, 30, pp. 434-439

PGM - Ion Semiconductor Sequencing

