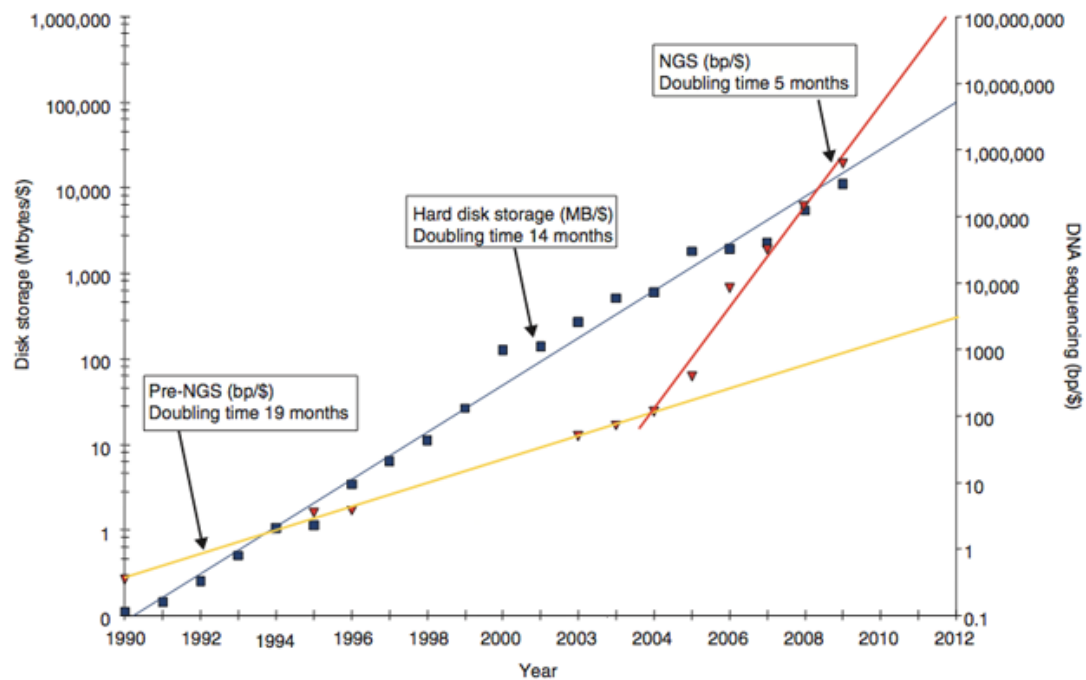


# Computational Genomics

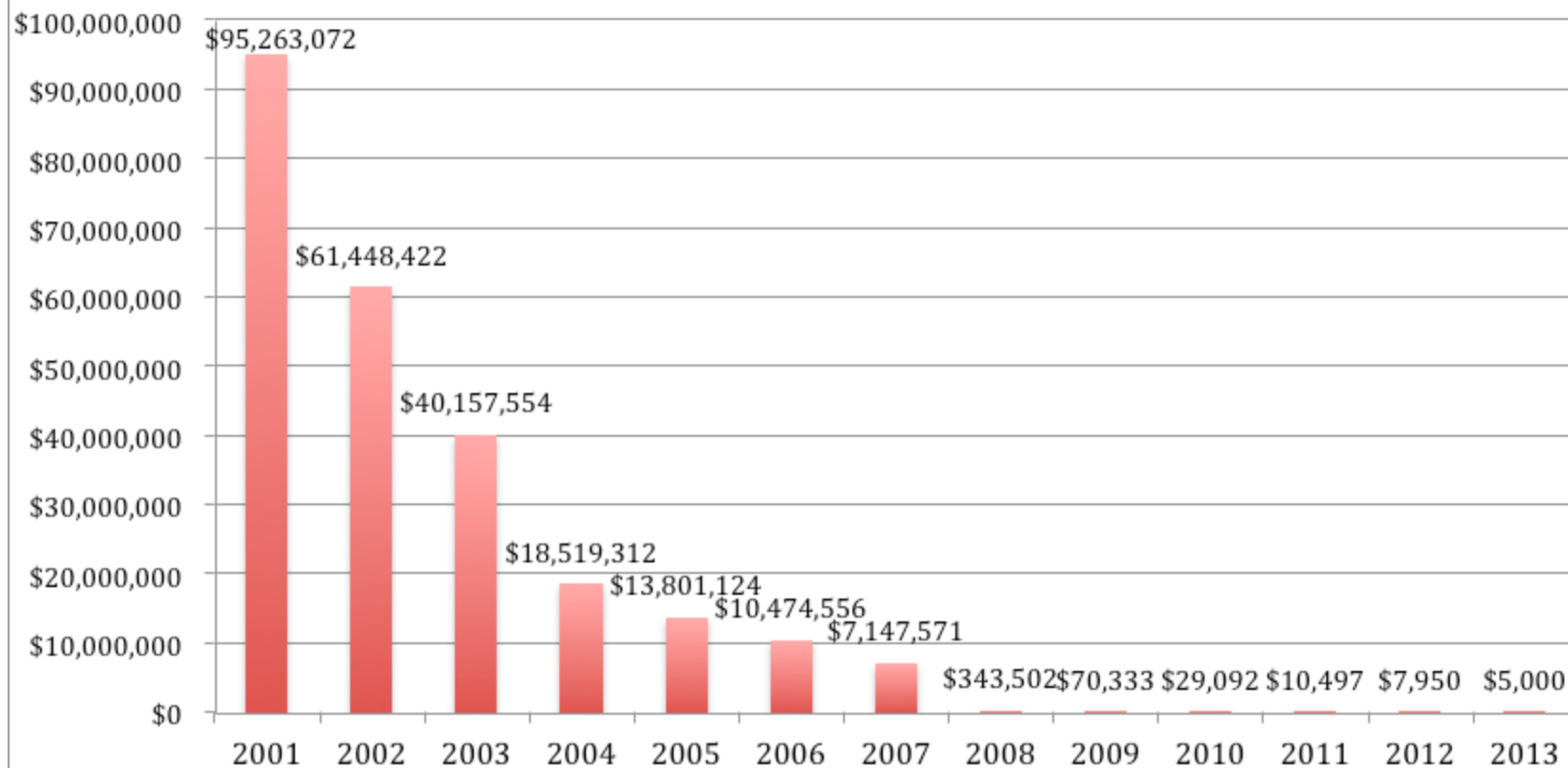
**Next generation sequencing (NGS)**

# Sequencing technology defies Moore's law

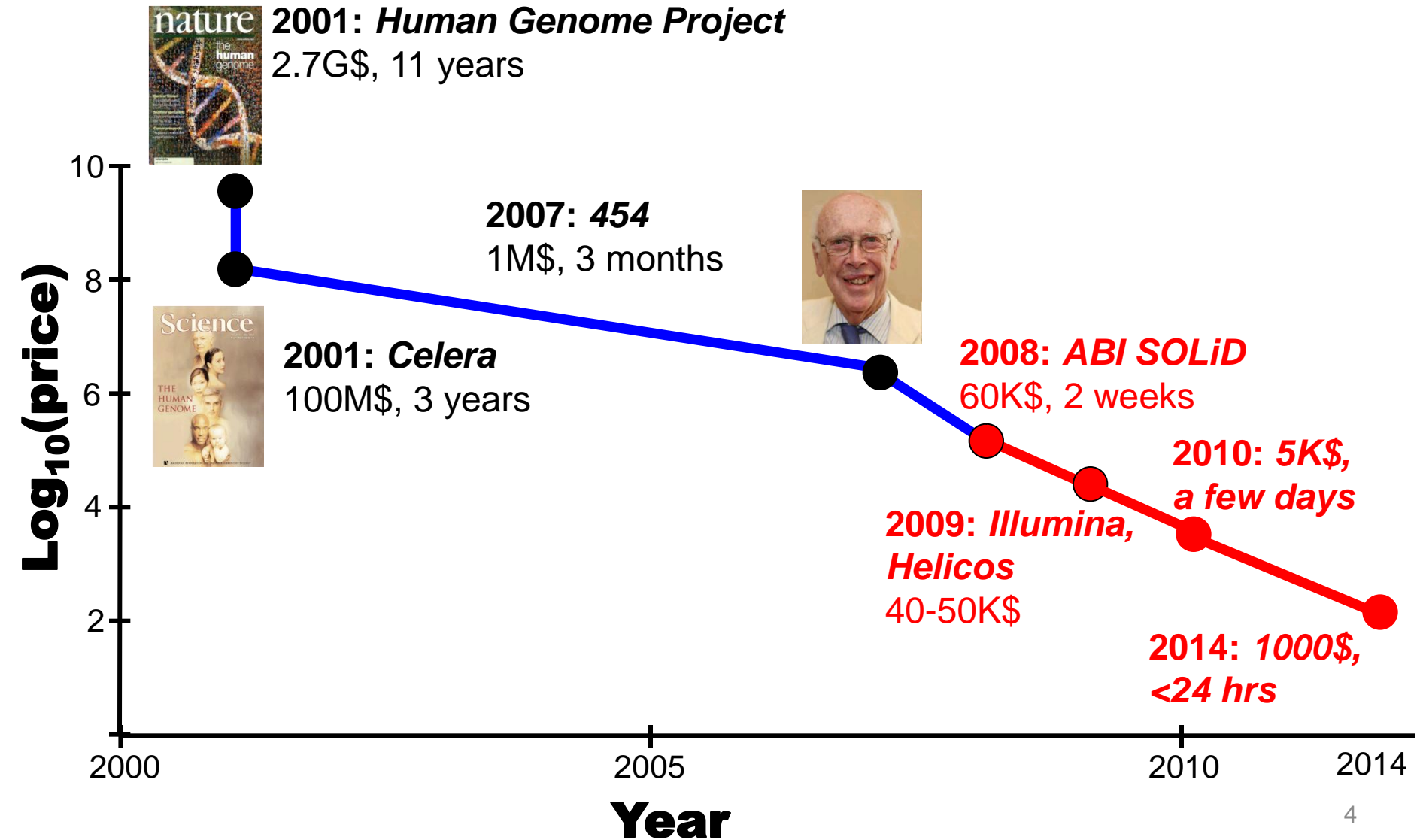


Nature Methods 2011

## Cost of Sequencing the Human Genome



# Sequencing the Human Genome





## Does Illumina Have the First \$1,000 Genome?

Illumina announces a new high-end sequencer made for “factory-scale” sequencing of human genomes.

By Susan Young on January 14, 2014



The \$1,000 genome has been a catchphrase of the sequencing industry for years, but despite bold promises from different companies, this benchmark hasn't been met. Now, thanks to a new sequencing machine from Illumina, it may finally be within reach.

At the J.P. Morgan Healthcare Conference on Tuesday, Illumina CEO Jay Flatley announced a new

### WHY IT MATTERS

Genome sequencing is still too costly for many medical applications, or for large-scale sequencing projects that could uncover the genetic basis of disease.

**Search Health** 3,000+ Topics

Go

**Inside Health**

[Research](#) | [Fitness & Nutrition](#)

## Company Unveils DNA Sequencing Device Meant to Be Portable, Disposable and Cheap

By [ANDREW POLLACK](#)

Published: February 17, 2012


DNA sequencing is becoming both faster and cheaper. Now, it is also becoming tinier.


A British company said on Friday that by the end of the year it would begin selling a disposable gene sequencing device that is the size of a USB memory stick and plugs into a laptop computer to deliver its

 [RECOMMEND](#)

 [TWITTER](#)

 [LINKEDIN](#)

 [SIGN IN TO E-MAIL](#)

 [PRINT](#)



# Next Gen-Omics™

**Genome  
Resequencing**

**mRNA Tag  
Profiling**

**Methylation  
Analysis**

**Small RNA  
Identification**

**Functional  
Elements**  
(ChIP-Seq, DNase-Seq)

**Transcriptome  
Sequencing**



# Applications of next-generation sequencing

Category	Examples of applications
Complete genome resequencing	Comprehensive polymorphism and mutation discovery in individual human genomes
Reduced representation sequencing	Large-scale polymorphism discovery
Targeted genomic resequencing	Targeted polymorphism and mutation discovery
Paired end sequencing	Discovery of inherited and acquired structural variation
Metagenomic sequencing	Discovery of infectious and commensal flora
Transcriptome sequencing	Quantification of gene expression and alternative splicing; transcript annotation; discovery of transcribed SNPs or somatic mutations
Small RNA sequencing	microRNA profiling
Sequencing of bisulfite-treated DNA	Determining patterns of cytosine methylation in genomic DNA
Chromatin immunoprecipitation–sequencing (ChIP-Seq)	Genome-wide mapping of protein-DNA interactions
Nuclease fragmentation and sequencing	Nucleosome positioning
Molecular barcoding	Multiplex sequencing of samples from multiple individuals

**Nature Biotechnology 26 (10): 1135-1145 (2008)**



# Roche & illumina analyzers

**Genome Sequencer FLX**



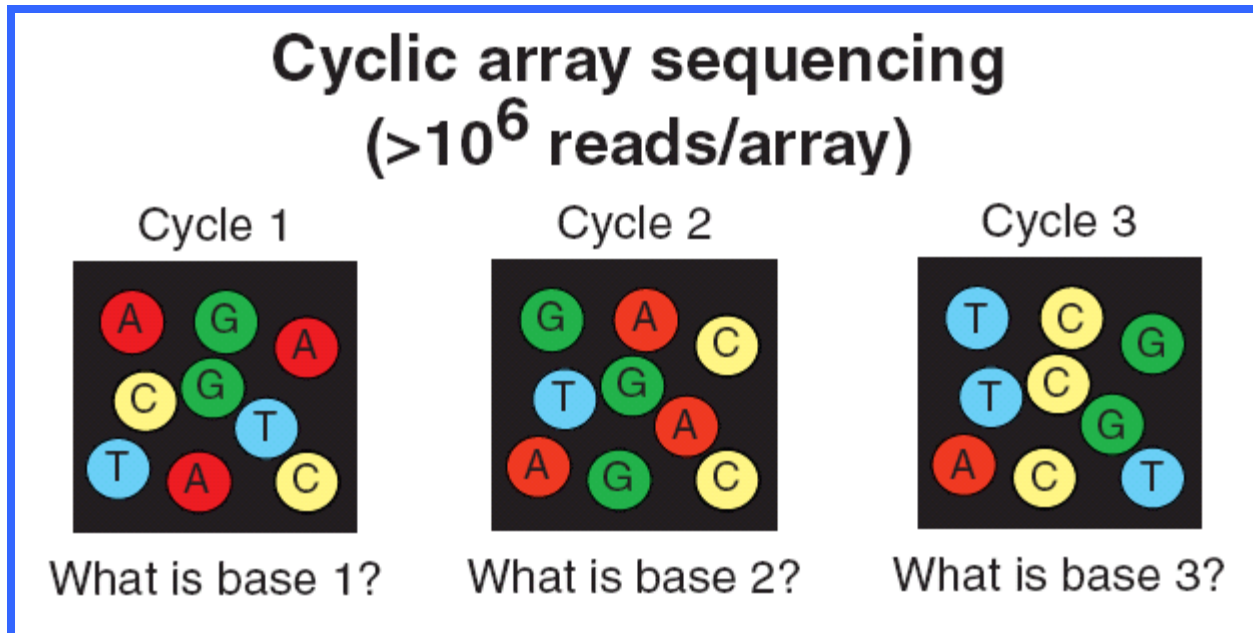
**Illumina Genome Analyzer**



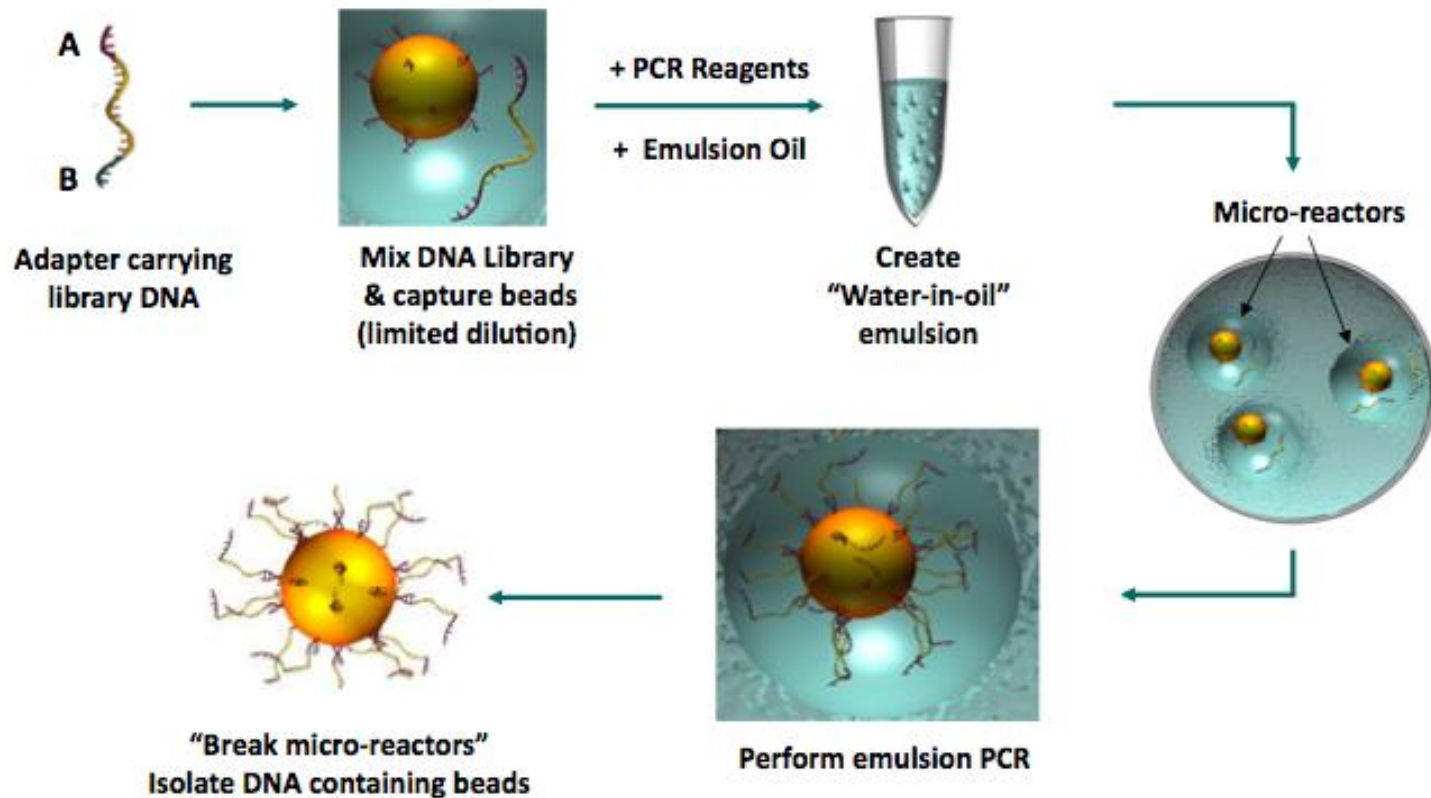
# High Parallelism is Achieved in Polony Sequencing

**Sanger**

**Polony**

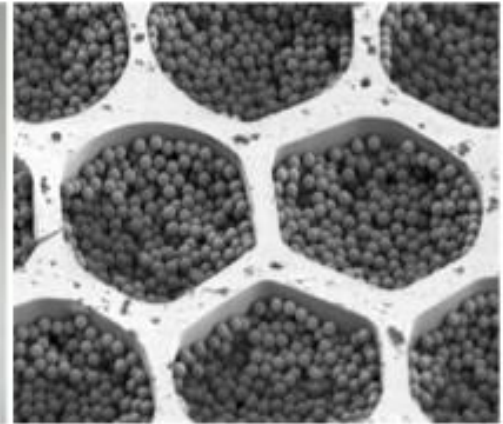
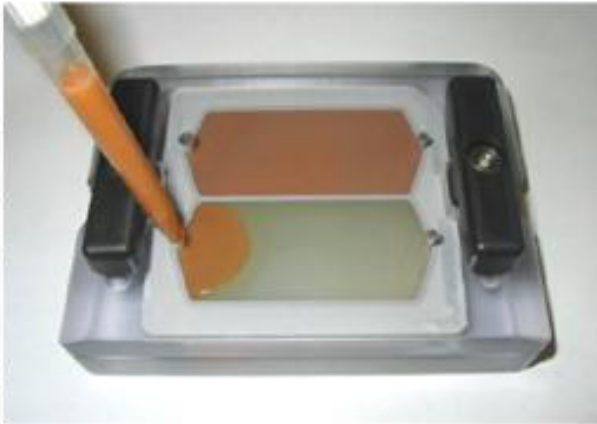
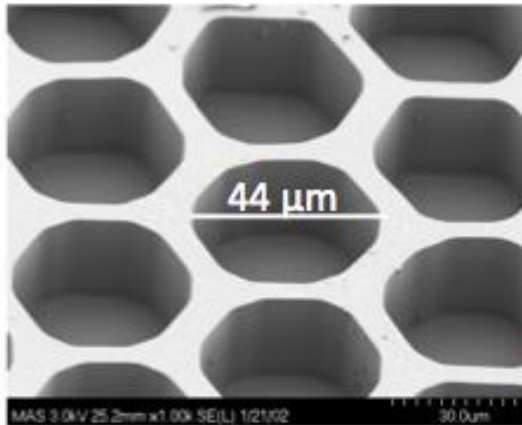
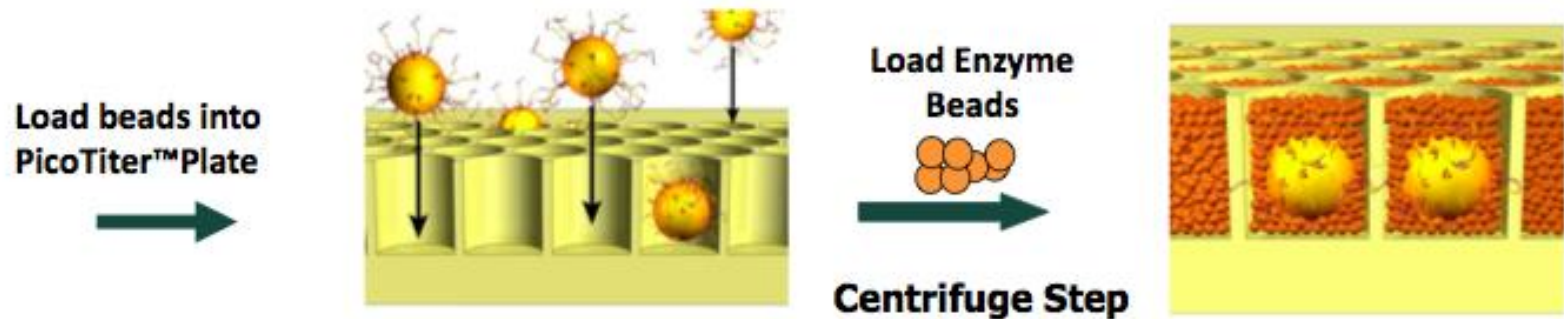


# Generation of Polony array: DNA Beads (454, SOLiD)



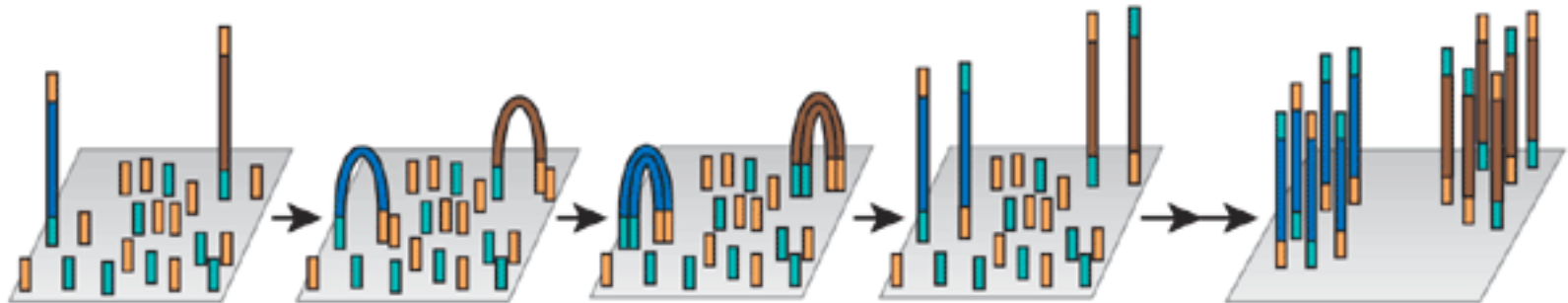
**DNA Beads are generated using Emulsion PCR**

# Generation of Polony array: DNA Beads (454, SOLiD)



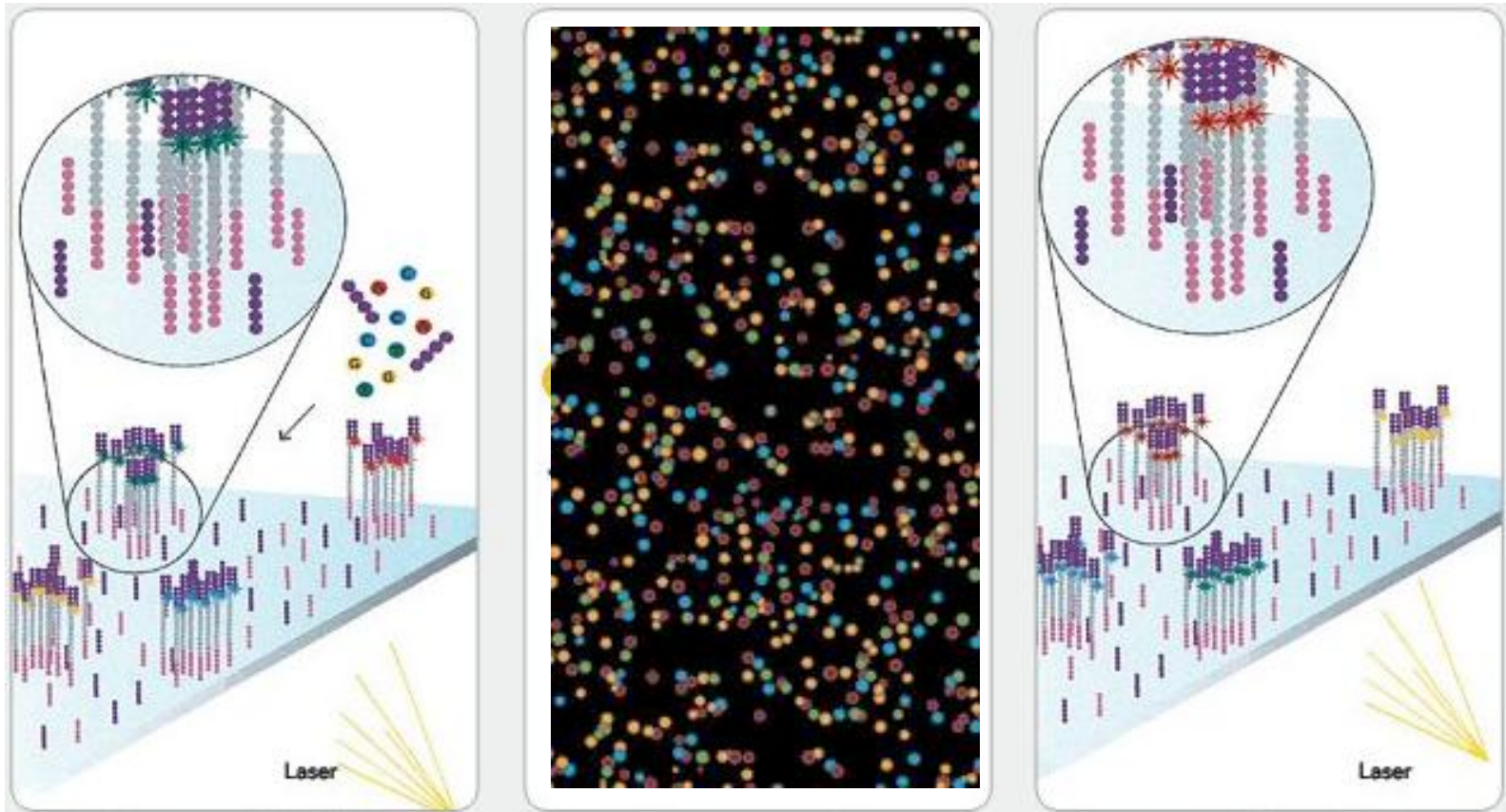
**DNA Beads are placed in wells**

# Generation of Polony array: Bridge-PCR (Solexa)



**DNA fragments are attached to array and  
used as PCR templates**

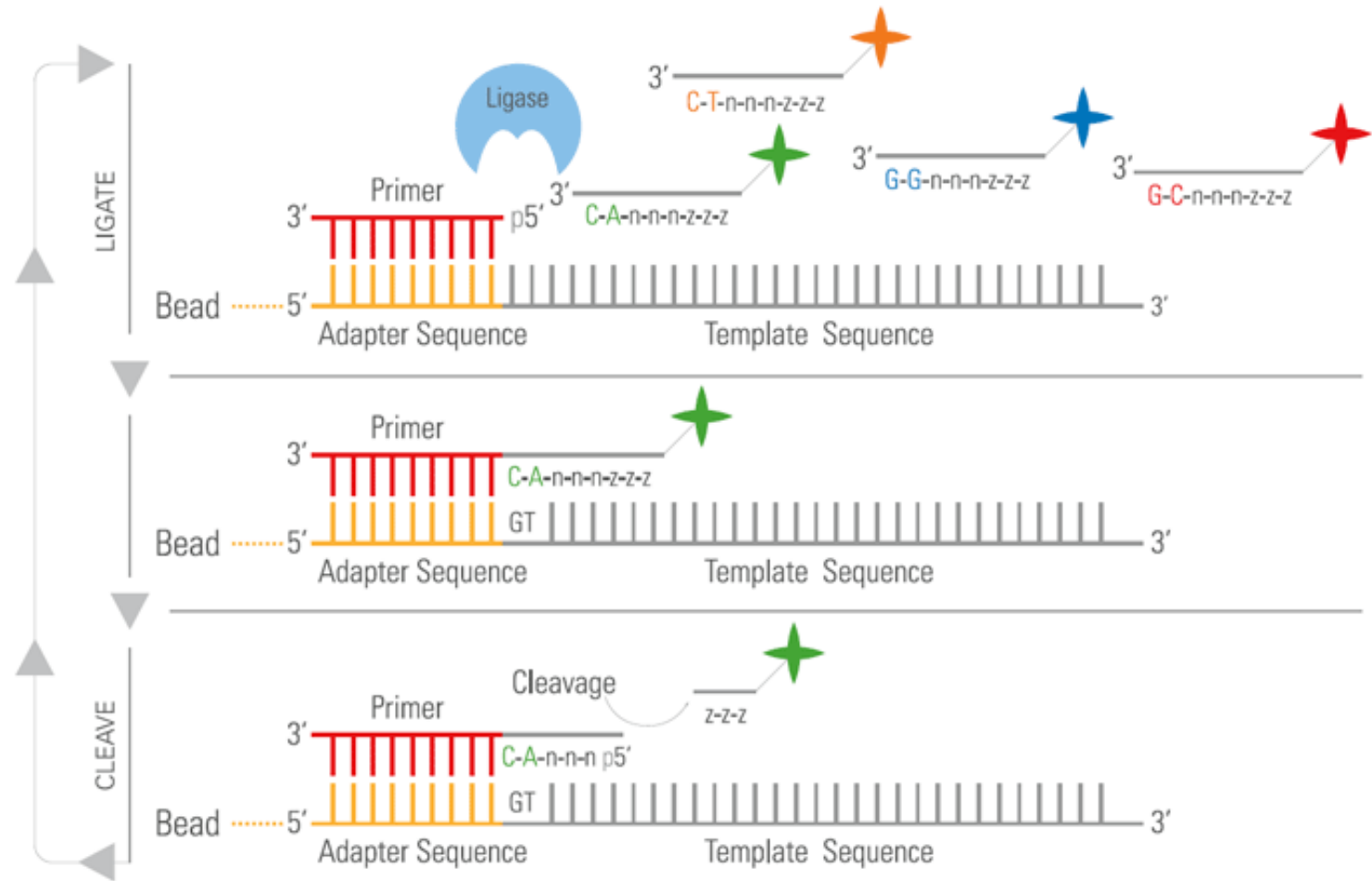
# Sequencing: Fluorescently labeled Nucleotides (Solexa)



**Complementary strand elongation: DNA Polymerase**



# Sequencing: Fluorescently Labeled Nucleotides (ABI SOLiD)



**Complementary strand elongation: DNA Ligase**

# **Single Molecule Sequencing: HeliScope**

- Direct sequencing of DNA molecules: no amplification stage
- DNA fragments are attached to array
- Potential benefits: higher throughput, less errors

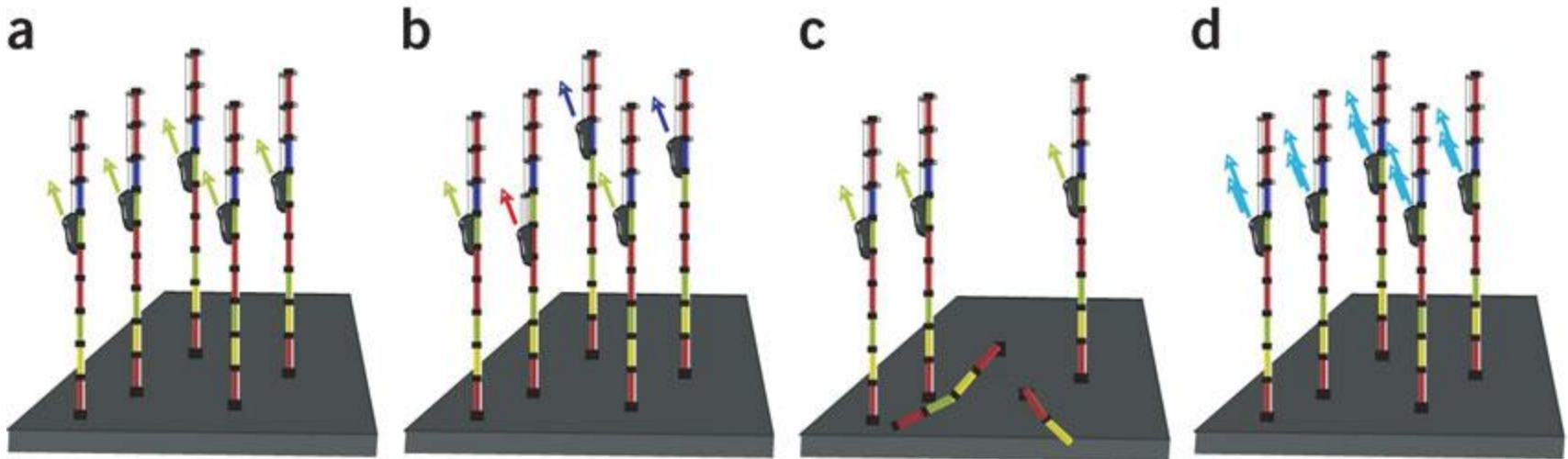


# Technology Summary

	Read length	Sequencing Technology	Throughput (per run)	Cost (1mbp)*
<b>Sanger</b>	~800bp	Sanger	400kbp	500\$
<b>454</b>	~400bp	Polony	500Mbp	60\$
<b>Solexa</b>	75bp	Polony	20Gbp	2\$
<b>SOLiD</b>	75bp	Polony	60Gbp	2\$
<b>Helicos</b>	30-35bp	Single molecule	25Gbp	1\$

\*Source: Shendure & Ji, *Nat Biotech*, 2008

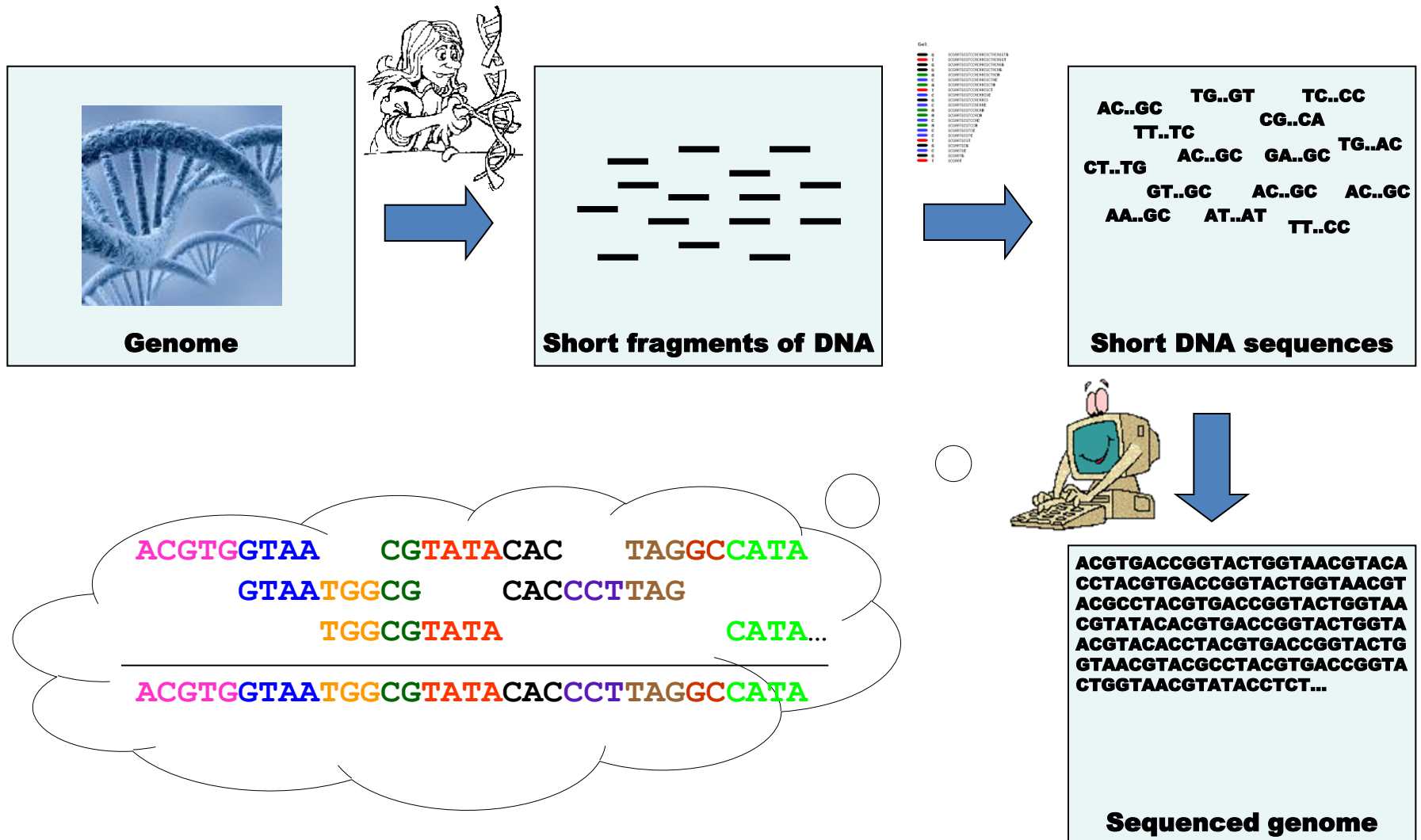
# Errors



# Genome Sequencing

- **Goal**
  - figuring the order of nucleotides across a genome
- **Problem**
  - Current DNA sequencing methods can handle only short stretches of DNA at once (usually between 100-200 bp's)
- **Solution**
  - Sequence and then use computers to assemble the small pieces

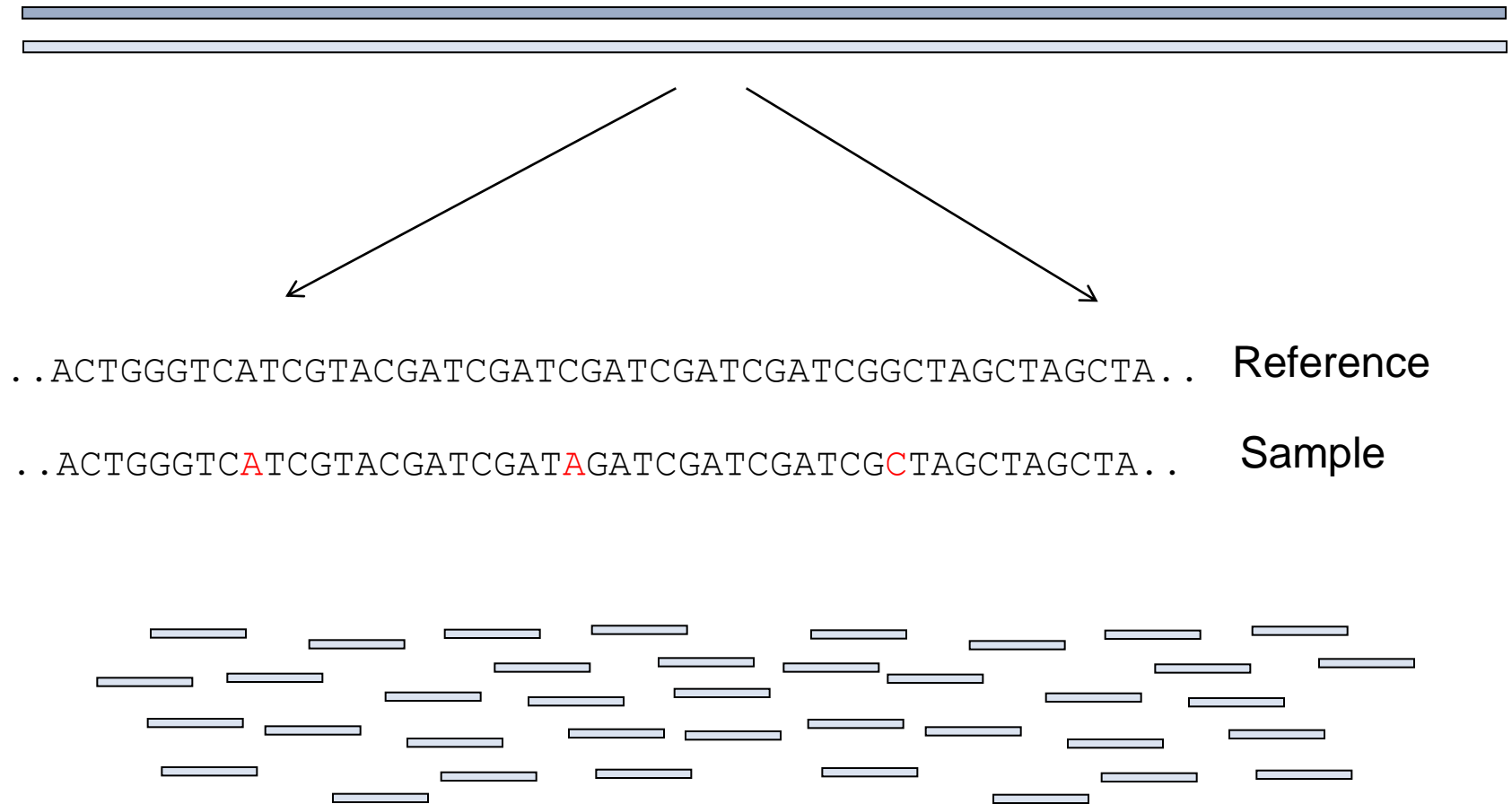
# Genome Sequencing



# Assembly

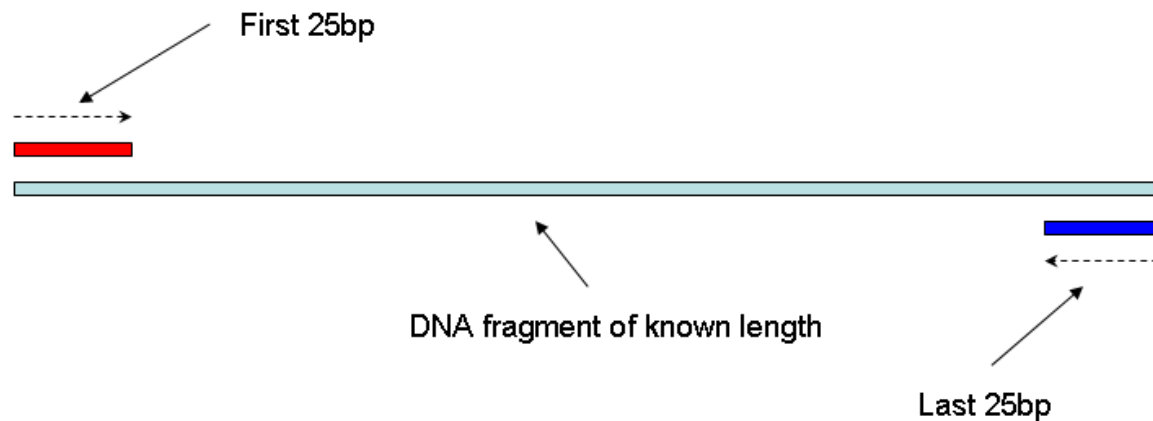
- How do we use the short reads to recover the genome being sequenced?

# 1: Using a reference genome: Alignment of reads to a reference



## 2. Assembly without a reference (de novo assembly)

- Paired-End sequencing (Mate pairs)
  - Sequence two ends of a fragment of known size.



- Currently fragment length (insert size) can range from 200 bps – 10,000 bps

# Velvet

- Euler / De Bruijn approach.
- Introduced as a alternative to overlap-layout-consensus approach in capillary sequencing.
- More suited for short read assembly.
- Based on De Bruijn graph.
- Implemented in *Velvet*<sup>1</sup>, the mostly used short read assembly method at present.

<sup>1</sup>Daniel Zerbino and Ewan Birney. Velvet: Algorithms for De Novo Short Read Assembly Using De Bruijn Graphs. Genome Res. 18: 821-829. 2008



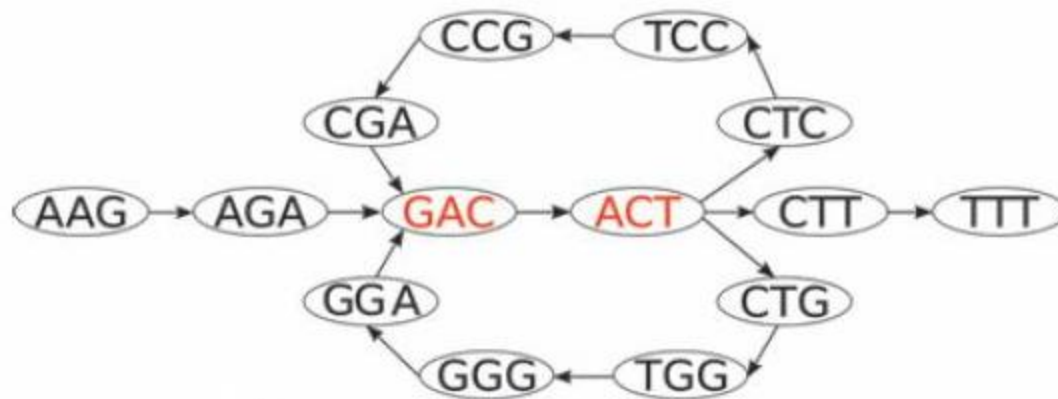
# De Bruijn graph method

- Break each read sequence in to overlapping fragments of size  $k$ . ( $k$ -mers)
- Form De Bruijn graph such that each  $(k-1)$ -mer represents a node in the graph.
- Edge exists between node  $a$  to  $b$  iff there exists a  $k$ -mer such that its prefix is  $a$  and suffix is  $b$ .
- Traverse the graph in unambiguous path to form contigs.

# De Bruijn graph

- $K = 4$

AA**GACT**CC**GACT**GG**GACT**TT



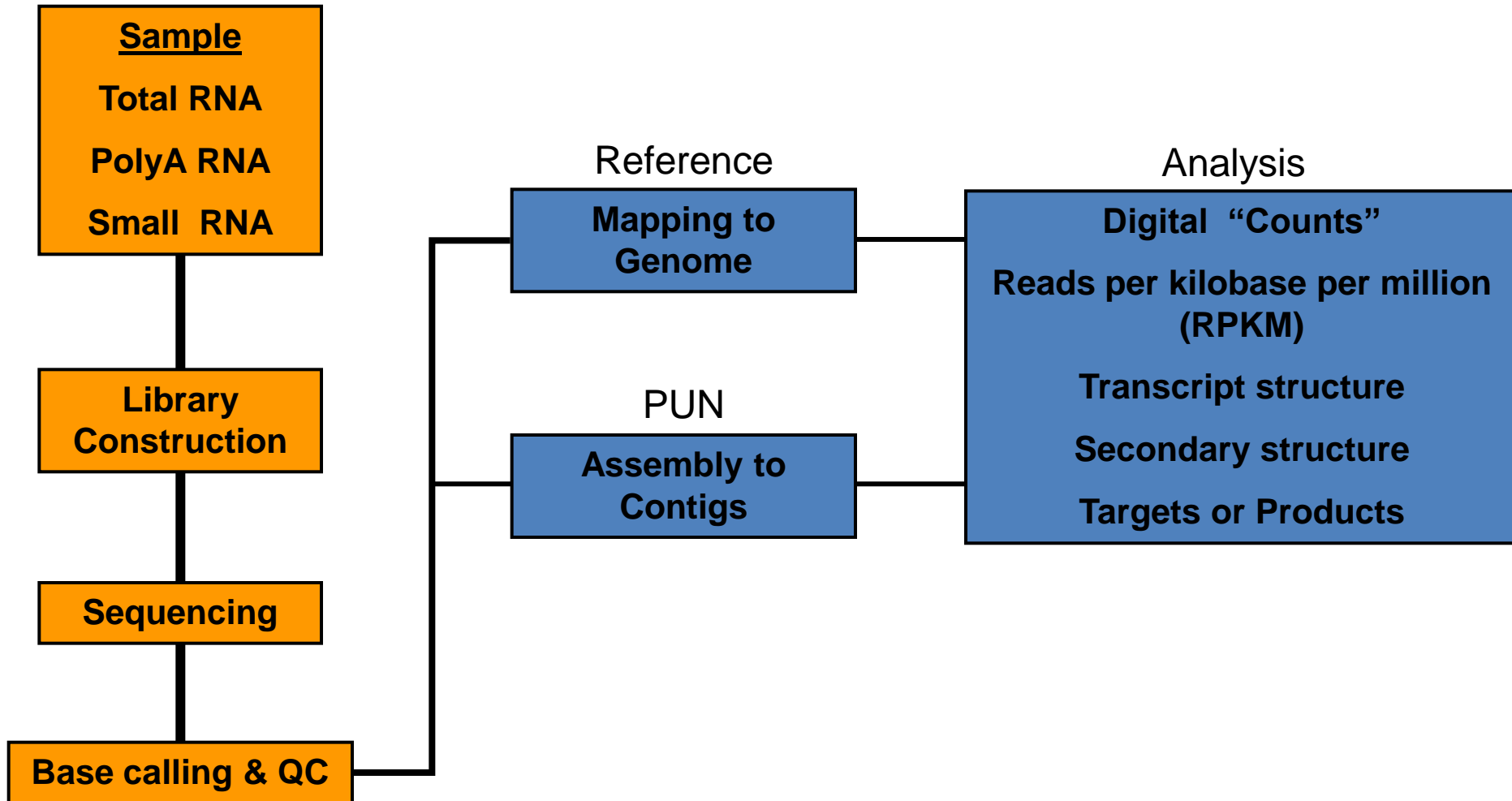
de Bruijn graph of a sequence

# De Bruijn graph method / Velvet

- Elegant way of representing the problem.
- Very fast execution.
- Error correction can be handled in the graph.
- De Bruijn graph size can be huge.
  - ~200GB for human genomes.
- Does not use pair information in initial phase, resulting in overly complicated graphs

Applications (beyond DNA)

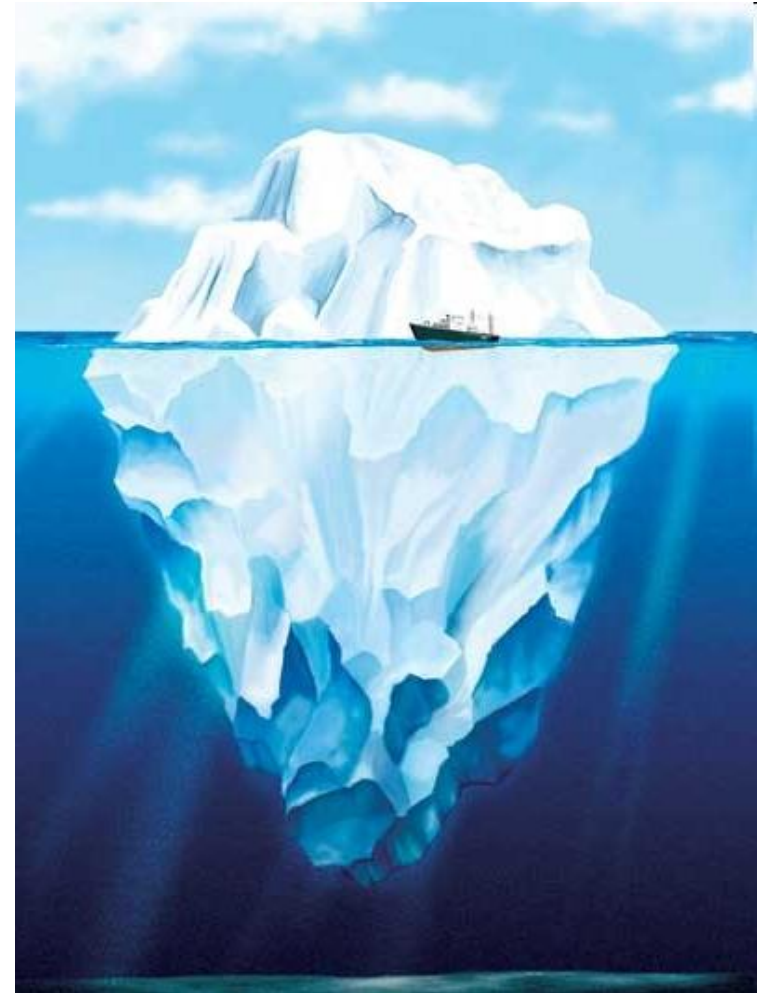
# RNASeq



# RNASeq – Compositional properties

## Depth of Sequence

- Sequence count  $\approx$  Transcript Abundance
- Majority of the data can be dominated by a small number of highly abundant transcripts
- Ability to observe transcripts of smaller abundance is dependent upon [sequence depth](#)



# Normalization

- Normalization is the process in which components of experiments are made comparable before statistical analysis.
- It is important in sequencing as well as for microarrays.
- A couple issues in normalization are different sequencing depth (library size) and distributions of reads (long right tails).

# Summarized Counts

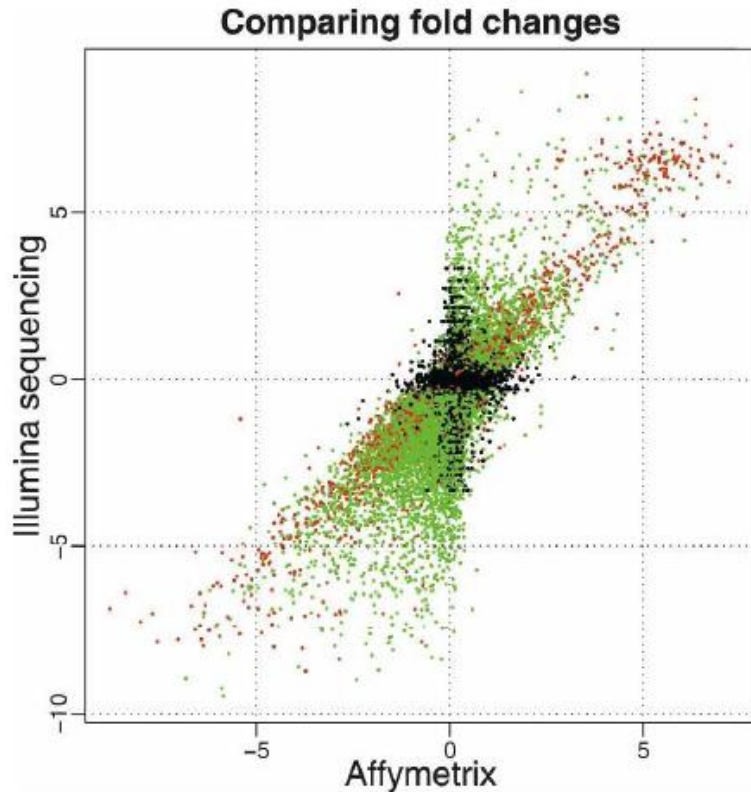
Gene	ORF size	Sample1	Sample2	Sample3	Sample4	Sample5	Sample6
<i>NDUFS2</i>	1861	1286	1700	1404	1485	1409	1161
<i>NDUFS3</i>	883	446	788	530	894	498	642
<i>NDUFS6</i>	417	521	745	722	843	651	609
<i>NFKBIE</i>	1549	59	68	79	52	74	30
<i>NOTCH4</i>	6157	4	2	2	2	5	0
<i>NRTN</i>	971	16	39	22	29	15	21
<i>YBX1</i>	1152	606	626	581	377	578	285
<i>ROR2</i>	3037	0	0	0	0	0	0
<i>OVOL1</i>	1126	0	0	0	0	0	0
<i>PARK2</i>	1541	2	0	0	0	1	0
<i>PCK2</i>	2061	102	171	97	108	89	87
<i>PET112L</i>	1721	184	316	197	280	215	265
<i>PEX14</i>	1161	194	211	259	179	189	115
<i>PFKFB3</i>	1953	155	84	280	109	240	51
<i>PFKFB4</i>	1433	151	79	297	106	213	68
<i>PIGH</i>	670	128	158	159	237	164	92
<i>PIK3C2G</i>	4463	0	0	1	0	0	0
<i>PKNOX1</i>	1517	87	141	109	219	113	156
<i>PKP2</i>	2767	154	150	114	128	122	95
<i>PLCB2</i>	3874	1	3	0	1	1	2
<i>SEPT4</i>	1686	1	6	1	6	0	8
<i>POU4F2</i>	1484	0	0	1	0	0	0
<i>PSPH</i>	1431	323	329	338	438	357	380
<i>RAB4A</i>	871	325	241	332	346	364	258
<i>MAP4K2</i>	2561	218	319	256	298	228	181



# Simple RPKM Normalization

- Proportion of reads: number of reads ( $n$ ) mapping to an exon (gene) divided by the total number of reads ( $N$ ),  $n/N$ .
- RPKM: Reads Per Kilobase of exon (gene) per Million mapped sequence reads,  $10^9 n / (NL)$ , where  $L$  is the length of the transcriptional unit in bp (Mortazavi et al., *Nat. Meth.*, 2008).

# RNASeq - Correspondence



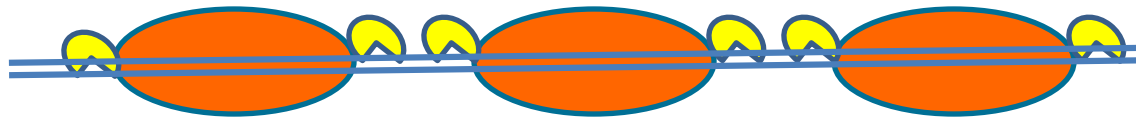
- Good correspondence with :
  - **Expression Arrays**
  - **Tiling Arrays**
  - **qRT-PCR**
- Range of up to 5 orders of magnitude
- Better detection of low abundance transcripts
- Greater power to detect
  - **Transcript sequence polymorphism**
  - **Novel trans-splicing**
  - **Paralogous genes**
  - **Individual cell type expression**

**Table 3 Bioinformatics tools for short-read sequencing**

Program	Categories	Author(s)	Reference	URL
Cross_match	Alignment	Phil Green, Brent Ewing and David Gordon		<a href="http://www.phrap.org/phredphrapconsed.html">http://www.phrap.org/phredphrapconsed.html</a>
ELAND	Alignment	Anthony J. Cox		<a href="http://www.illumina.com/">http://www.illumina.com/</a>
Exonerate	Alignment	Guy S. Slater and Ewan Birney	72	<a href="http://www.ebi.ac.uk/~guy/exonerate">http://www.ebi.ac.uk/~guy/exonerate</a>
MAQ	Alignment and variant detection	Heng Li	37	<a href="http://maq.sourceforge.net">http://maq.sourceforge.net</a>
Mosaik	Alignment	Michael Strömberg and Gabor Marth		<a href="http://bioinformatics.bc.edu/marthlab/Mosaik">http://bioinformatics.bc.edu/marthlab/Mosaik</a>
RMAP	Alignment	Andrew Smith, Zhenyu Xuan and Michael Zhang	73	<a href="http://rulai.cshl.edu/rmap">http://rulai.cshl.edu/rmap</a>
SHRIMP	Alignment	Michael Brudno and Stephen Rumble		<a href="http://compbio.cs.toronto.edu/shrimp">http://compbio.cs.toronto.edu/shrimp</a>
SOAP	Alignment	Ruiqiang Li <i>et al.</i>	35	<a href="http://soap.genomics.org.cn">http://soap.genomics.org.cn</a>
SSAHA2	Alignment	Zemin Ning <i>et al.</i>	36	<a href="http://www.sanger.ac.uk/Software/analysis/SSAHA2">http://www.sanger.ac.uk/Software/analysis/SSAHA2</a>
SXOligoSearch	Alignment	Synamatix		<a href="http://synasite.mgrc.com.my:8080/sxog/NewSXOligoSearch.php">http://synasite.mgrc.com.my:8080/sxog/NewSXOligoSearch.php</a>
ALLPATHS	Assembly	Jonathan Butler <i>et al.</i>	38	
Edena	Assembly	David Hernandez <i>et al.</i>	74	<a href="http://www.genomic.ch/edena">http://www.genomic.ch/edena</a>
Euler-SR	Assembly	Mark Chaisson and Pavel Pevzner	75	
SHARCGS	Assembly	Juliane Dohm <i>et al.</i>	76	<a href="http://sharcgs.molgen.mpg.de">http://sharcgs.molgen.mpg.de</a>
SHRAP	Assembly	Andreas Sundquist <i>et al.</i>	39	
SSAKE	Assembly	René Warren <i>et al.</i>	40	<a href="http://www.bcgsc.ca/platform/bioinfo/software/ssake">http://www.bcgsc.ca/platform/bioinfo/software/ssake</a>
VCAKE	Assembly	William Jeck	77	<a href="http://sourceforge.net/projects/vcake">http://sourceforge.net/projects/vcake</a>
Velvet	Assembly	Daniel Zerbino and Ewan Birney	41	<a href="http://www.ebi.ac.uk/%7Ezerbino/velvet">http://www.ebi.ac.uk/%7Ezerbino/velvet</a>
PyroBayes	Base caller	Aaron Quinlan <i>et al.</i>	34	<a href="http://bioinformatics.bc.edu/marthlab/PyroBayes">http://bioinformatics.bc.edu/marthlab/PyroBayes</a>
PbShort	Variant detection	Gabor Marth		<a href="http://bioinformatics.bc.edu/marthlab/PbShort">http://bioinformatics.bc.edu/marthlab/PbShort</a>
ssahaSNP	Variant detection	Zemin Ning <i>et al.</i>		<a href="http://www.sanger.ac.uk/Software/analysis/ssahaSNP">http://www.sanger.ac.uk/Software/analysis/ssahaSNP</a>

Incomplete list compiled from sources, including <http://seqanswers.com/forums/showthread.php?t=43> and <http://www.sanger.ac.uk/Users/lh3/seq-nt.html>.

# ChIPSeq



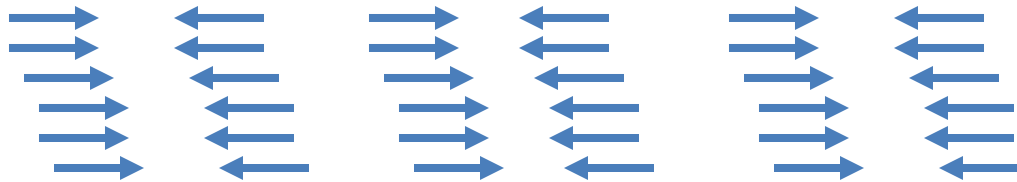
MNase  
Linker Digest



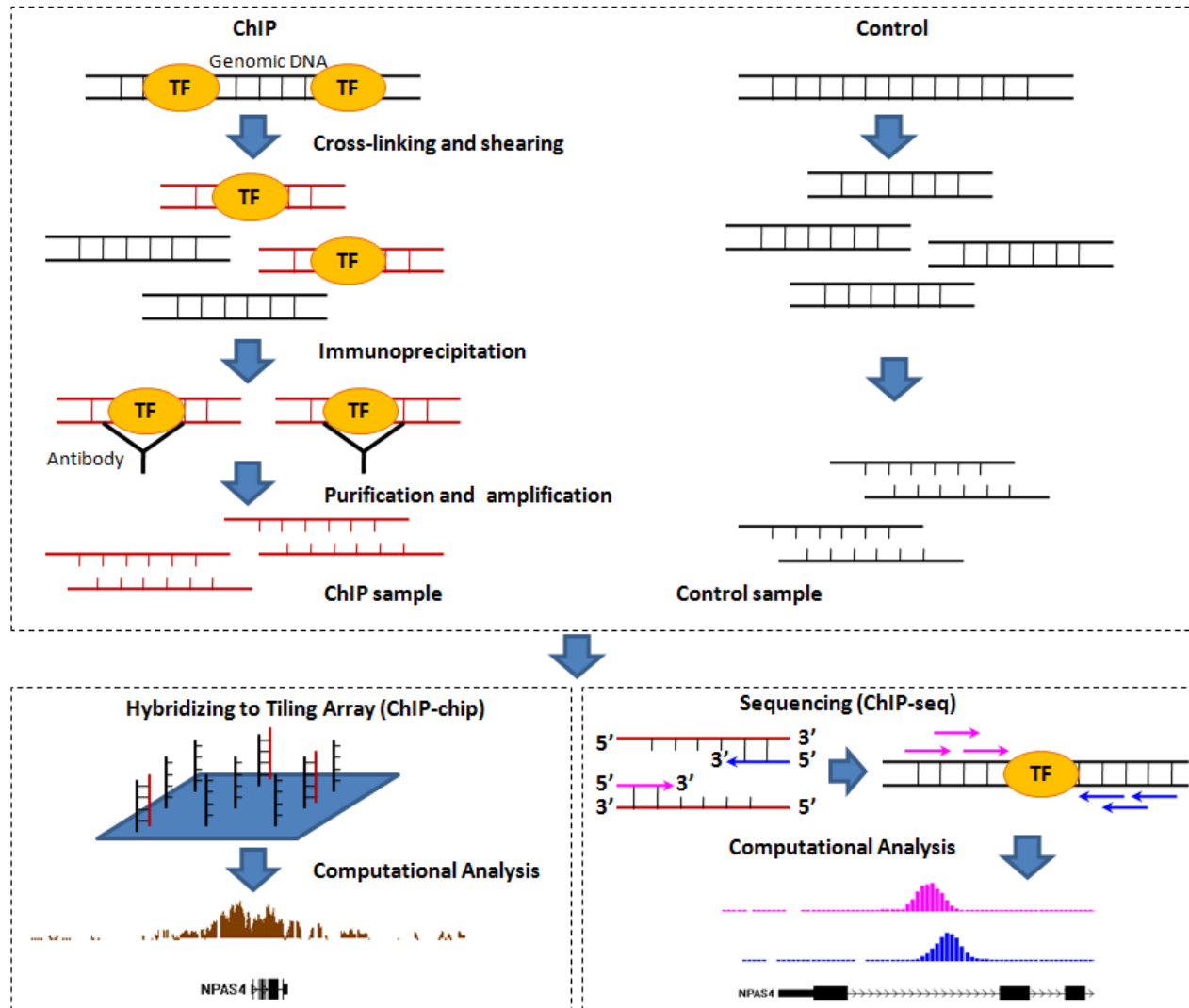
Remove  
Nucleosomes



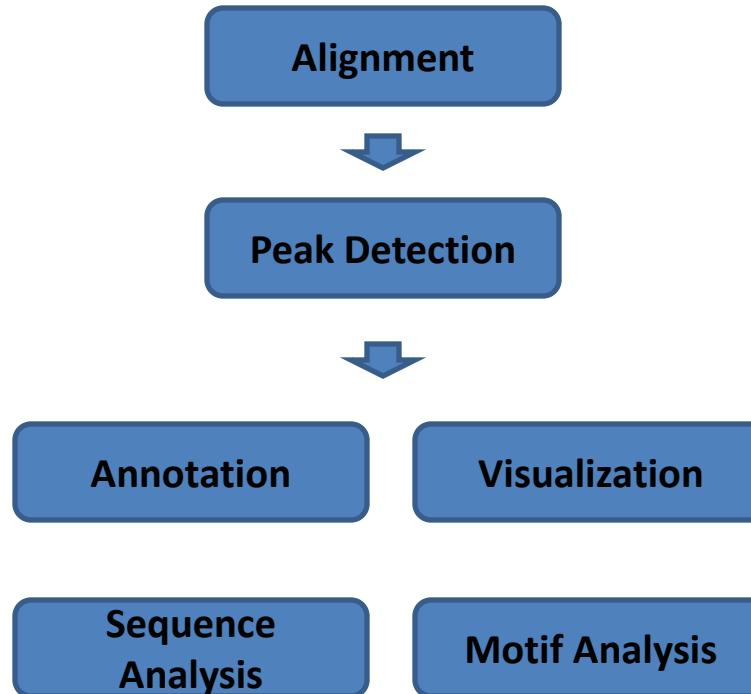
Sequence &  
Align



# ChIP-Seq



# ChIP-Seq Analysis



# Alignment

- ELAND
- Bowtie
- SOAP
- SeqMap
- ...

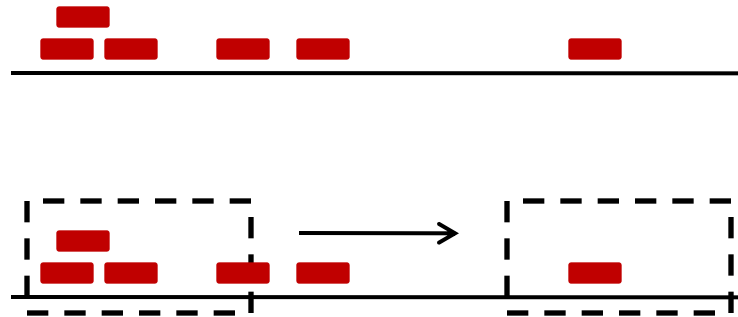
# Peak detection

- FindPeaks
- CHiPSeq
- BS-Seq
- SISSRs
- QuEST
- MACS
- CisGenome
- ...



# One sample analysis

A simple way is the sliding window method



Poisson background model is commonly used to estimate error rate

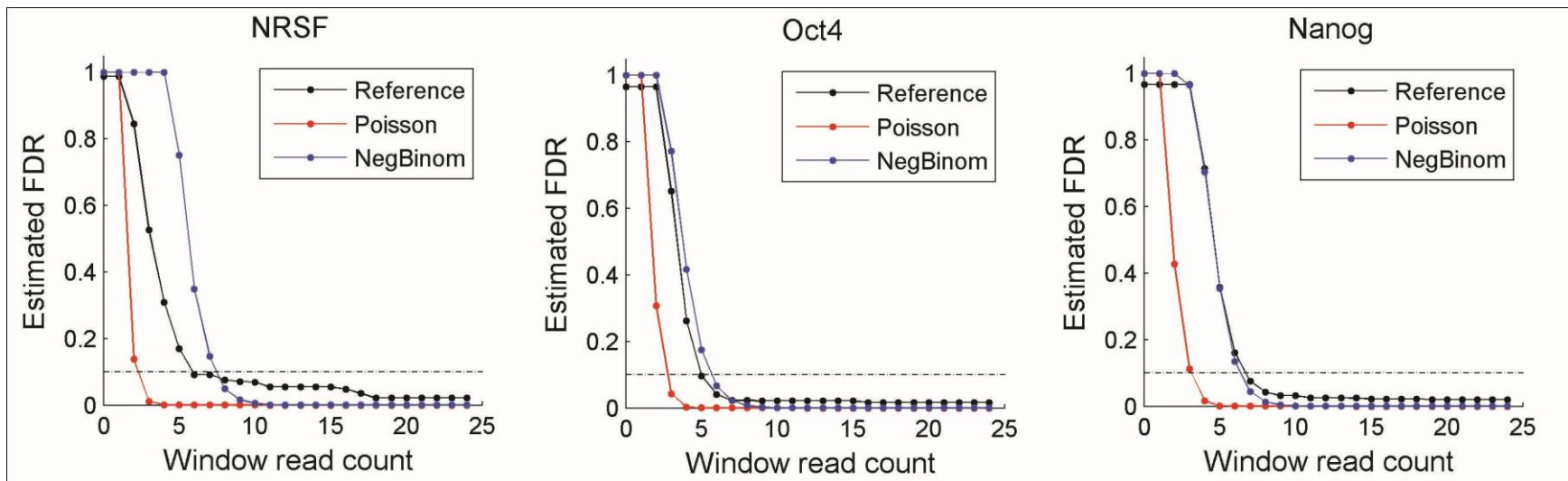
$$k_i \sim \text{Poisson}(\lambda_0)$$



Or people use Monte Carlo simulations

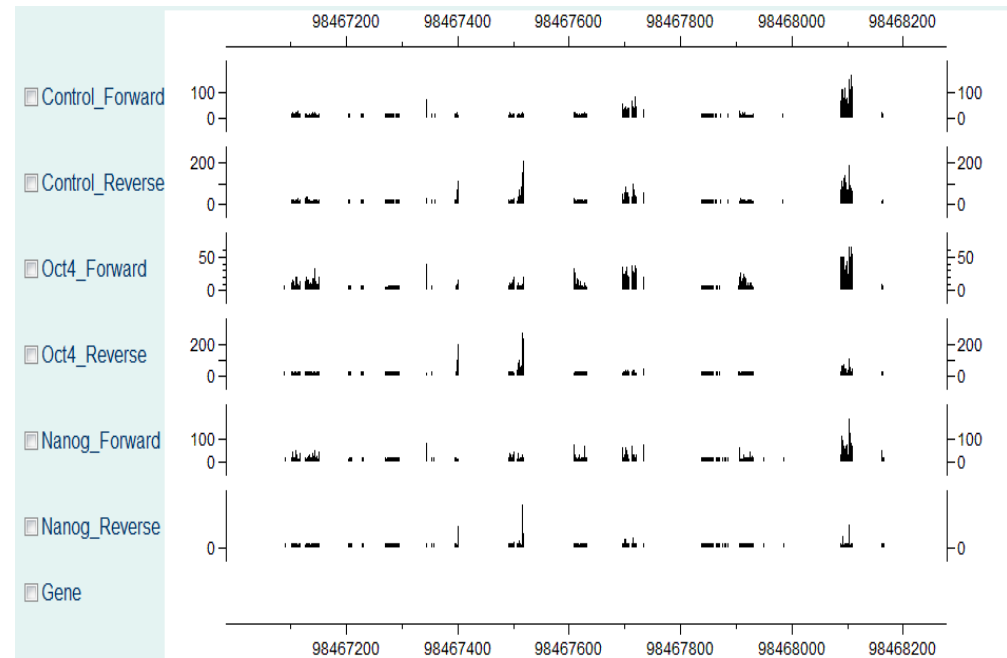
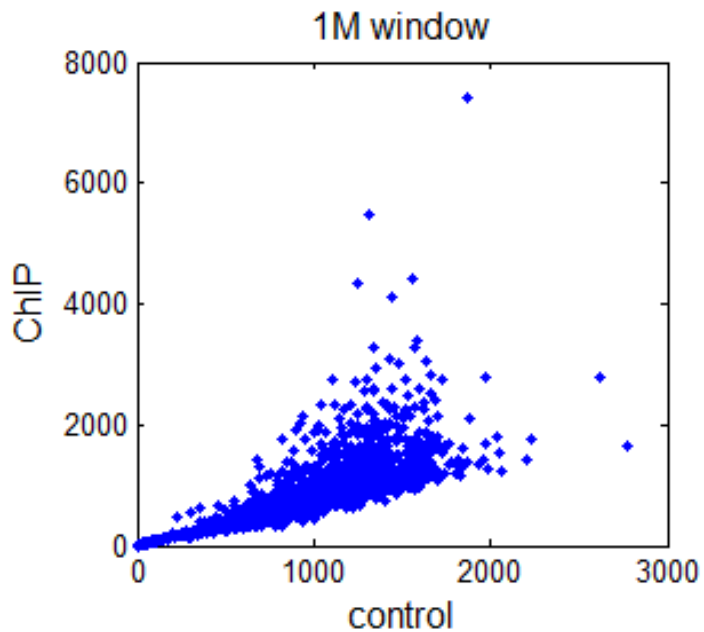
Both are based on the assumption that read sampling rate is a constant across the genome.

# FDR estimation based on Poisson and negative binomial model

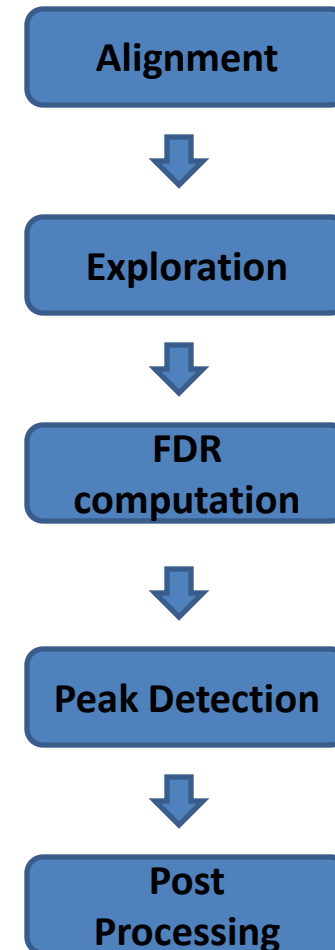
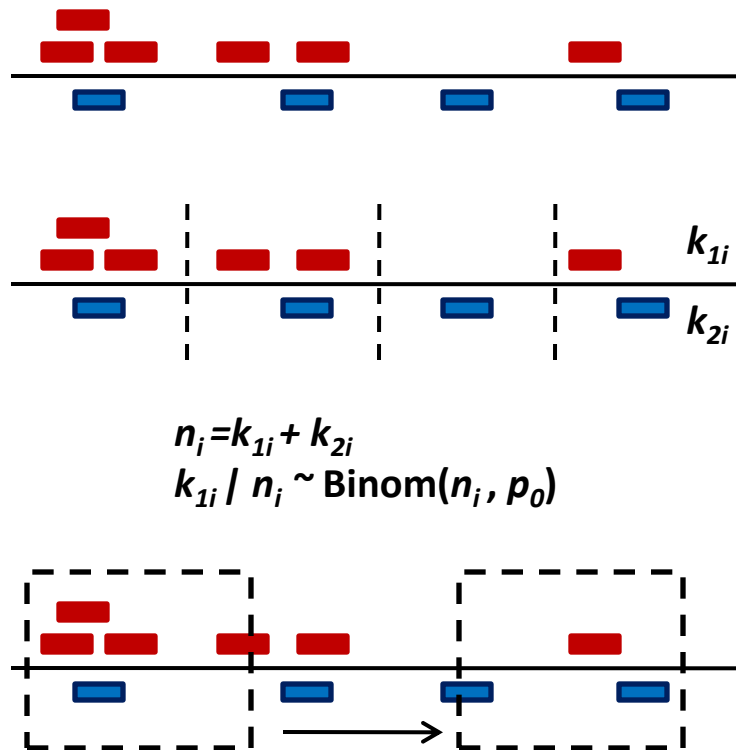


# Two sample analysis

Reason: read sample rates at the same genomic locus are correlated across different samples.



# CisGenome two sample analysis



# Epigenome

- Protein-DNA interactions [ChIPSeq]
  - Nucleosome positioning
  - Histone modification
  - Transcription factor interactions
- Methylation [MethylSeq]
- Impact of NextGen
  - Whole genome profiling
  - Resolution

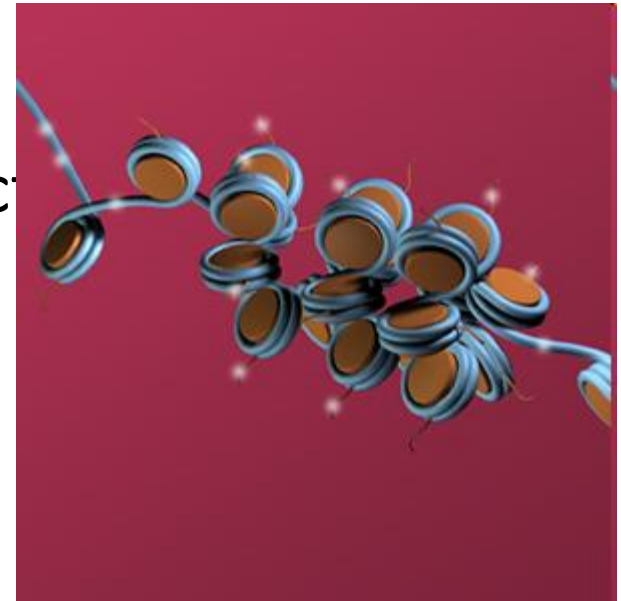


Image : ClearScience

# Third generation sequencing

## 3<sup>rd</sup> generation

- Single Molecule Sequencing Technology (tSMS)
- No amplification
- 10Gb os sequence data per 8 day run
- Single Molecule Real Time (SMRT) sequencing technology (PacBio RS)

## PacBio RS



# Nanopore sequencing

