

# Introduction to NGS Technologies

**Ignacio Medina**

**im411@cam.ac.uk**

**Head of Computational Biology Lab  
HPC Service, University of Cambridge, UK**

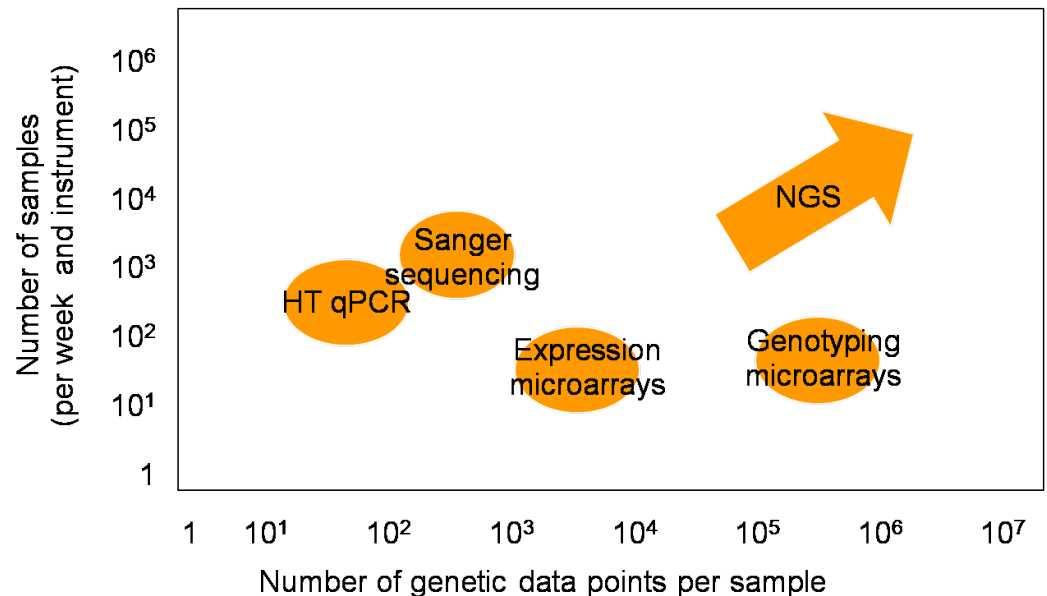
**EMBL-EBI Scientific collaborator  
Genome Campus, Hinxton, Cambridge, UK**



**UNIVERSITY OF  
CAMBRIDGE**

# Relative throughput of the different HT technologies

NGS emerges with a potential of data production that will, eventually wipe out conventional HT technologies in the years coming



Too many sequences to be handled in a standard computer

	<b>Sanger (1st-gen) Sequencing</b>	<b>Next-Gen Sequencing, and 3rd generation</b>
Whole Genome	Human (early drafts), model organisms, bacteria, viruses and mitochondria (chloroplast), low coverage	New human (!), individual genome, exomes, 2,500 normal (1K genome project), 25,000 cancer (TCGA and ICGC initiatives), CNV, matched control pairs, time course, rare-samples
RNA	cDNA clones, ESTs, Full Length Insert cDNAs, other RNAs	RNA-Seq: Digitization of transcriptome, alternative splicing events, miRNA, allele specific transcripts
Communities	Environmental sampling, 16S RNA populations, ocean sampling,	Human microbiome, deep environmental sequencing, Bar-Seq
Other		Epigenome, rearrangements, ChIP-Seq

# NGS technologies



**Cost-effective  
Fast  
Ultra throughput  
Cloning-free  
Short reads**



# Differences between the various platforms:

- Nanotechnology used.
- Read length
- Chemistry and enzymology.
- Signal to noise detection in the software
- Software/images/file size/pipeline
- Cost
- Applications

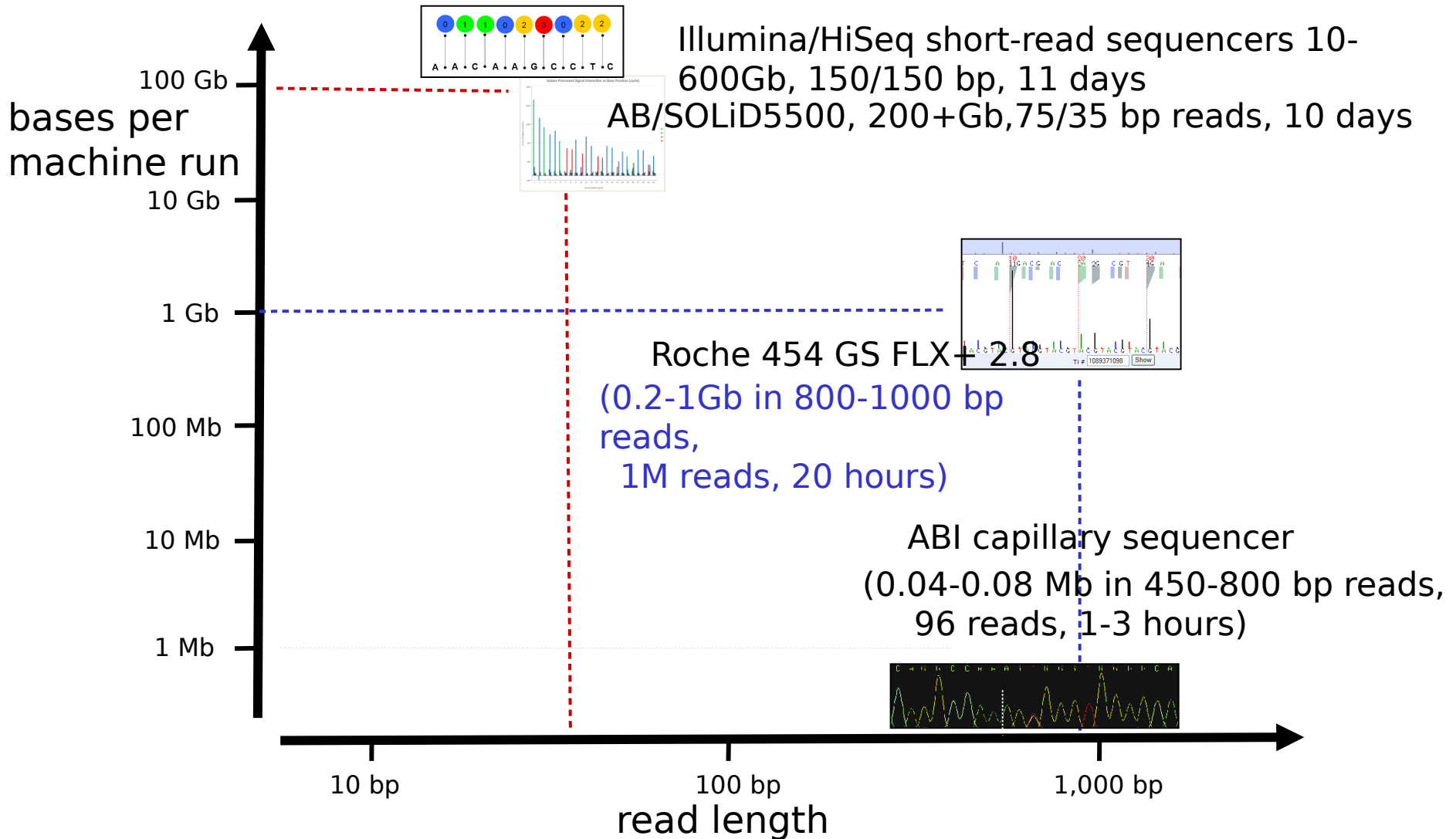
# Similarities- LOTS of DATA

## General ways of dealing at the sequences

- Assemble them and look at what you have
- You map them (align against a known genome) and then look at what you have.
- Or a mixture of both!
- Sometimes you select the DNA you are sequencing
- or you try to sequence everything
- Depends on biological question, sequencing machine you have, and how much time and money you have.
- **NGS is relatively cheap but think what you want to answer, because the analysis won't do magic**

# Next-gen sequencers

From John McPherson, OICR



# Next Generation Sequencers

In the past 3 main platforms:

- **Solexa/illumina**
- **Roche 454**
- **ABI SOLiD**
- Follow an approach similar to Sanger sequencing, but do away with separation of fragments by size and “read” the sequence as the reaction occurs
- Several different “next generation” sequencing platforms developed and commercialized, more on the way.
- Simultaneously sequence entire libraries of DNA sequence fragments



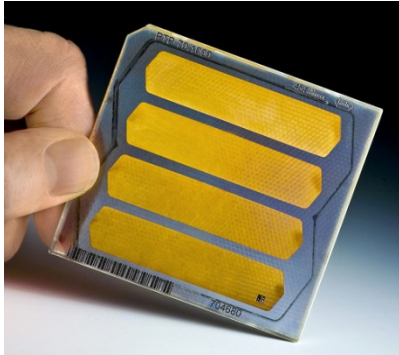
# 454 (Roche)

- First next generation method to be commercially available
- Uses a “sequencing by synthesis” (SBS) approach:
  - DNA is broken into pieces of 500-1,400 bp, ligated to adaptors, and amplified on tiny beads by PCR (emulsion PCR)
  - Beads (with DNA attached) are placed into tiny wells (one bead per well) on a PicoTiter Plate that has millions of wells. Each well is connected to an optical fibre.
  - DNA is sequenced by adding polymerase and DNA bases containing pyrophosphate. The different bases (A,C,G,T) are added sequentially in a flow chamber
  - When a base complementary to the template is added, the pyrophosphate is released and a burst of light is produced
  - The light is detected and used to call the base
- Initially 100-150 bp, but they have been improved to 600-1000 bp
- >1 million, filter-passed reads per run (20 hours)
- 1 billion bases per day

# Roche 454: GS FLX System

- Good for
  - “*de novo*” sequencing (longer reads)
  - Resequencing (expensive)
  - New bacterial genomes.
  - Amplicons
- Pyrosequencing. Bias with long polynucleotide stretches

# Roche 454 GS FLX



<b>Throughput</b>	400-600 million high-quality, filter-passed bases per run* 1 billion bases per day
<b>Run Time</b>	10 hours
<b>Read Length</b>	Average length = 400 bases
<b>Accuracy</b>	Q20 read length of 400 bases (99% at 400 bases and higher for prior bases)
<b>Reads per run</b>	>1 million high-quality reads
<b>Data</b>	Trace data accepted by NCBI since 2005
<b>Computing Requirements</b>	Cluster recommended (Roche GS FLX Titanium Cluster available)
<b>Robustness</b>	No complex optics or lasers; reagents have long shelf life



# GS Junior, benchtop



## System Performance

<b>Throughput</b>	35 million high-quality, filtered bases per run*
<b>Run Time</b>	10 hours sequencing 2 hours data processing
<b>Avg. Read Length</b>	400 bases*
<b>Accuracy</b>	Q20 read length of 400 bases (99% accuracy at 400 bases)
<b>Reads per Run</b>	100,000 shotgun, 70,000 amplicon
<b>Sample Input</b>	gDNA, amplicons, cDNA, or BACs depending on the application
<b>Physical Dimensions</b>	40 cm wide x 60 cm deep x 40 cm high (the size of a laser printer) Weight = 55 lbs.
<b>Computing</b>	Linux-based OS on HP desktop computer included. All software is point-and-click.

*\*Typical results. Average read length and number of reads depend on specific sample and genomic characteristics*

# Illumina

- Over 90% of all sequencing data is produced on Illumina systems.
- Uses a “sequencing by synthesis” approach:
  - DNA is broken into small fragments and ligated to an adaptor.
  - The fragments are attached to the surface of a flow cell and amplified.
  - DNA is sequenced by adding polymerase and labeled reversible terminator nucleotides (each base with a different color).
  - The incorporated base is determined by fluorescence.
  - The fluorescent label is removed from the terminator and the 3' OH is unblocked, allowing a new base to be incorporated
- Started with 35 bp, increased now to up to 150 bp
- One run can give up to 10-600 Gb, 300-6000 million paired-end reads
- 75-85% of bases at or above Q30

## Illumina HiSeq 2500



600 Gb/run in 11 days  
2x100 bp fragments  
3-6 billion reads per run

# Illumina Systems

## Illumina MiSeq



**175-245 Mb** 4h 1x 36bp  
**1.5-2.0 Gb** 27h 2x150 bp

## Illumina HiSeq X Ten

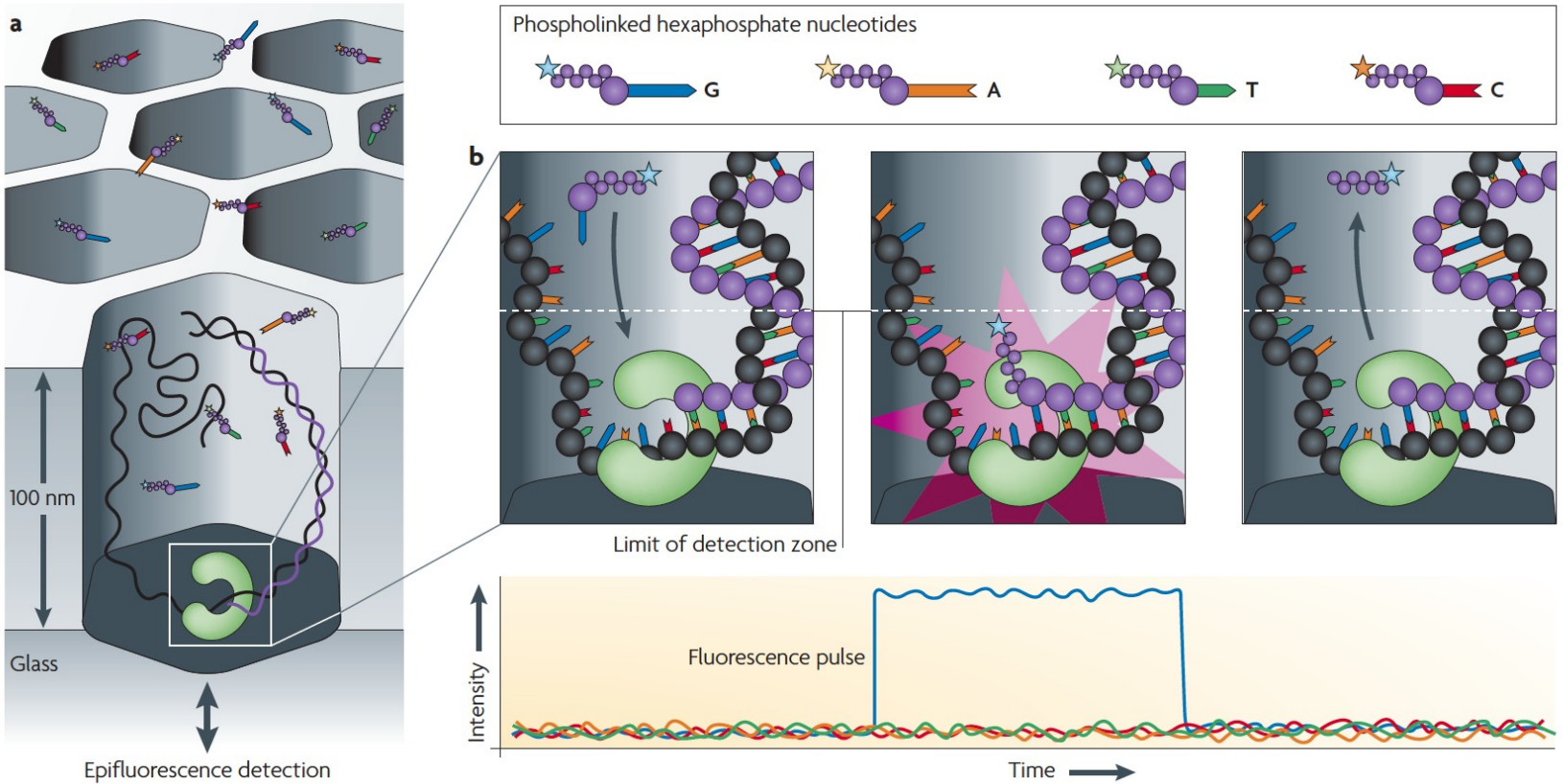


1800 Gb/run in <3 days  
2x150 bp fragments  
6 billion reads per run

Consists of a set of **10 HiSeq X** ultra-high-throughput instruments that deliver over **18,000 human genomes** per year at the price of **\$1000** per genome.

# PacBio

## Pacific Biosciences — Real-time sequencing



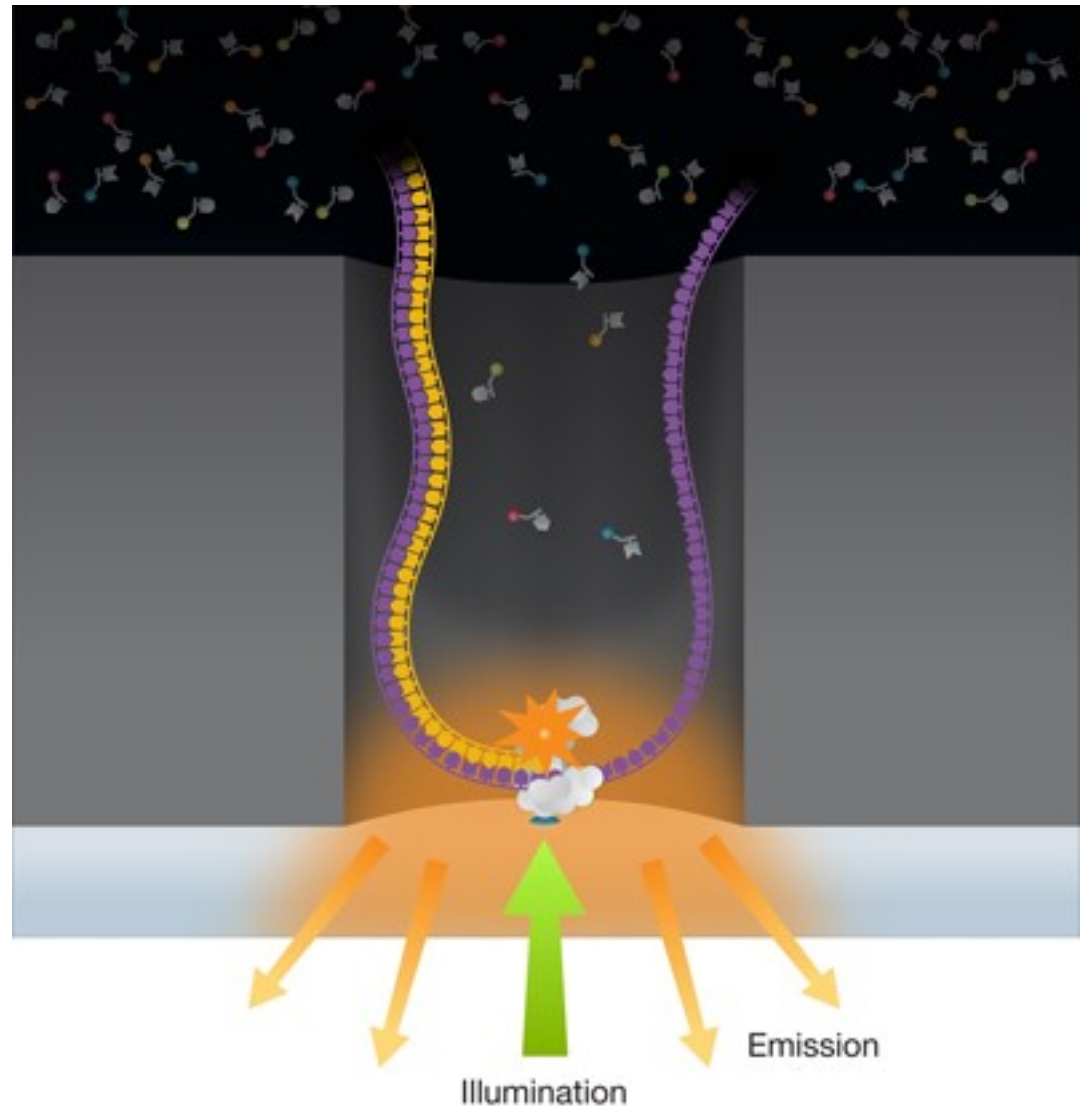
From Michael Metzker, <http://view.ncbi.nlm.nih.gov/pubmed/19997069>

# Pacific Bioscience

SMRT: Singel Molecule Real  
time DNA synthesis

Up to 12000 nt  
50 bases/second

ZMW: Zero Mode Waveguide





# Ion Torrent

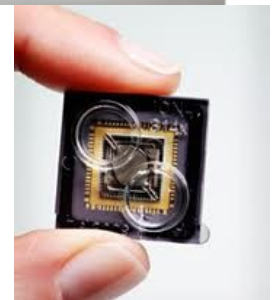
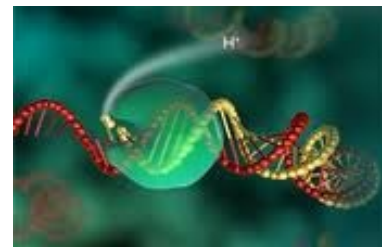
\$ 50.000

\$ 500 /sample

1 hour/run

> 200 nt lengths

Reads H<sup>+</sup> released by DNA  
polymerase



# Comparison

## Roche 454

- Long fragments
- Errors: poly nts
- Low throughput
- Expensive
  
- De novo sequencing
- Amplicon sequencing
- RNASeq

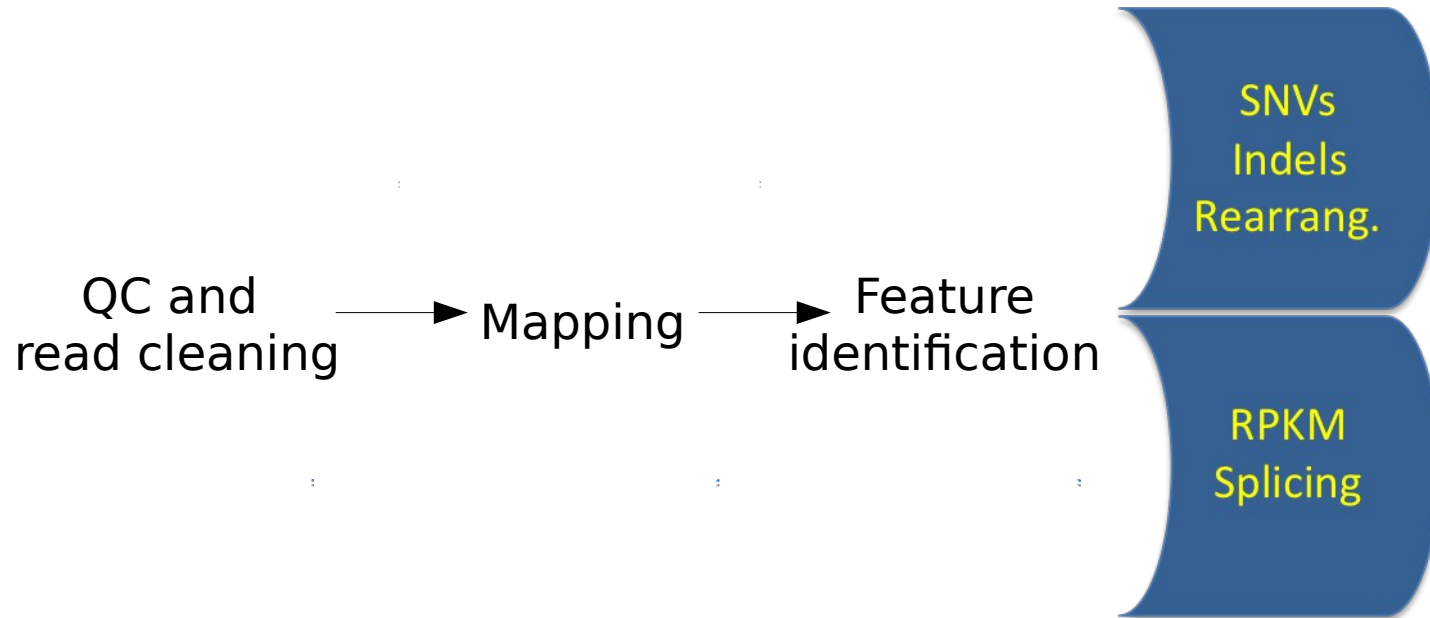
## Illumina

- Short fragments
- Errors: Hexamer bias
- High throughput
- Cheap
  
- Resequencing
- De novo sequencing
- ChipSeq
- RNASeq
- MethylSeq

## SOLiD

- Short fragments
- Color-space
- High throughput
- Cheap
  
- Resequencing
- ChipSeq
- RNASeq
- MethylSeq

# Basic steps NGS data processing



# File formats

```
+ILLUMINA-AAATTTTATTTTAACTTGTAAGAAGGTGTCGTT
+ILLUMINA-GA_0000:1:1:4010:1065#0/1
hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh
@ILLUMINA-GA_0000:1:1:4093:1065#0/1
AAATAACTAAGAAATTTGTCACAAATTTCTCAAATTCCTT
+ILLUMINA-GA_0000:1:1:4093:1065#0/1
!ffffffggggaaffccfdffcdffdgffgcggfvgggg
jcarbonell@ender:/scratch2/jcarbonell$
jcarbonell@ender:/scratch2/jcarbonell$ head -n 20
@ILLUMINA-GA_0000:1:1:1395:1061#0/2
GGCACAAGCAAGCAAGTGTCTGAATTCCTTTGCAGAGATA
+ILLUMINA-GA_0000:1:1:1395:1061#0/2
hcaehghcne_WffffiffafajffcgghghgeahheWfff
@ILLUMINA-GA_0000:1:1:1855:1066#0/2
GTGTAATTCCTGTGCGCGGTTTTATGTGATGCCATCCA
+ILLUMINA-GA_0000:1:1:1855:1066#0/2
ffffcffffdhdfcfddffjjcc```dfffcchha
@ILLUMINA-GA_0000:1:1:3567:1062#0/2
TGTAGTCGGCGCGGAGCAAGCTGCCAGCCCCACCOCCECA
+ILLUMINA-GA_0000:1:1:3567:1062#0/2
hhhhhhhhhhhhhgfcfcffjddffS!efffhhhhhh'
@ILLUMINA-GA_0000:1:1:4010:1065#0/2
TTTGTTTGACAGTTAATGATGGTCTATTACATAAACGT
+ILLUMINA-GA_0000:1:1:4010:1065#0/2
hhhhhhghghghhghghghfhhghhhhhhhhhhhhhfhfe
@ILLUMINA-GA_0000:1:1:4093:1065#0/2
AATCCACAAGAGCAAAACAGTTGCCAAGAGATGCAAGGAC
+ILLUMINA-GA_0000:1:1:4093:1065#0/2
dfffffffhdhhhhghgfhhchggh_fQfbfffffidfa
jcarbonell@ender:/scratch2/jcarbonell$
jcarbonell@ender:/scratch2/jcarbonell$ samtools view ivial5_06_pair1.remdup
ILLUMINA-GA_0000:1:1:1395:1061#0    99      scaffold_13     799896   0
M:i:1 X0:i:0 XG:i:0 MD:Z:6A31
ILLUMINA-GA_0000:1:1:1395:1061#0    147     scaffold_13     800074   0
147 XM:i:1 X0:i:0 XG:i:0 MD:Z:2LC16
ILLUMINA-GA_0000:1:1:1855:1066#0    89      scaffold_65     576129   0
7 XM:i:2 X0:i:0 XG:i:0 MD:Z:3G4A29
ILLUMINA-GA_0000:1:1:3567:1062#0    83      scaffold_215    8768     0
M:i:1 X0:i:0 XG:i:0 MD:Z:3RL6
ILLUMINA-GA_0000:1:1:3567:1062#0    163     scaffold_215    8554     0
62 XM:i:2 X0:i:0 XG:i:0 MD:Z:18T1GL7
ILLUMINA-GA_0000:1:1:4010:1065#0    99      scaffold_76     865926   60
0 XM:i:0 X0:i:0 XG:i:0 MD:Z:38
ILLUMINA-GA_0000:1:1:4010:1065#0    147     scaffold_76     866076   60
0 XM:i:2 X0:i:0 XG:i:0 MD:Z:2C24AI0
ILLUMINA-GA_0000:1:1:4093:1065#0    99      scaffold_57     479190   12
2 XM:i:1 X0:i:0 XG:i:0 MD:Z:12G25
ILLUMINA-GA_0000:1:1:4093:1065#0    147     scaffold_57     479954   20
2 XM:i:0 X0:i:0 XG:i:0 MD:Z:38
ILLUMINA-GA_0000:1:1:6805:1068#0    99      scaffold_11     3541452  0
1 X0:i:0 XG:i:0 MD:Z:8A29
```

fastq: sequence data and qualities

## SAM/BAM: mapping data and qualities

```

carbone11@bender:~/scratch2/jcarbone11$ samtools view -h -l 1:1:1395:1061#0 99 scaffold_13 799896 0 38M = 800074 216 AATAGANACCACATTTGTAACAACTTTAGTCGCTGTTTC affffaBa``cc^ccfc_ffcfdfffc|ffddbbdfcc XT:A:R NM:i:1 SM:i:0 AM:i:0 X0:i:261 X
ILLUMINA-GA_0000:1:1:1395:1061#0 147 scaffold_13 800074 0 38M = 799896 -216 TATCTCTGCAAGAATTAGCATTTCTTGGTCTGTC ffwheahghgggfcff|afaffffffw_echgheach XT:A:R NM:i:1 SM:i:0 AM:i:0 X0:i:3 X1:i:
147 XM:i:1 X0:i:0 XG:i:0 MD:Z:21C16
ILLUMINA-GA_0000:1:1:1855:1066#0 89 scaffold_65 576129 0 38M = 576129 0 TTTTTTCTCTCTTTTGGCCATATTCTTCTCCTT cX|cfffaw`c`ccff|ffgggfffd|fd|fcffff XT:A:R NM:i:2 SM:i:0 AM:i:0 X0:i:2 X1:i:
7 XM:i:2 X0:i:0 XG:i:0 MD:Z:3G4A29
ILLUMINA-GA_0000:1:1:3567:1062#0 83 scaffold_215 8768 0 38M = 8554 -252 CCCCAGGCTATAGCCACCCGCTTTTTGGGNATTTT gfggggggfffffffcgggggeeeeeeeBeggggg XT:A:R NM:i:1 SM:i:0 AM:i:0 X0:i:250 X
M:i:1 X0:i:0 XG:i:0 MD:Z:31C6
ILLUMINA-GA_0000:1:1:3567:1062#0 163 scaffold_215 8554 0 38M = 8768 252 TGAATCGCGCGGACGACGTCGCGAGCCCCACCCCCCA hhhhhhhghhhghcgfcff|fdffS|effchhhhhh XT:A:R NM:i:2 SM:i:0 AM:i:0 X0:i:4 X1:i:
62 XM:i:2 X0:i:0 XG:i:0 MD:Z:18T1G17
ILLUMINA-GA_0000:1:1:4010:1065#0 99 scaffold_76 865926 60 38M = 866076 188 AAATAAAAAATATTTATTTAACTTCTAAGCATGTCGT hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh XT:A:U NM:i:0 SM:i:37 AM:i:37 X0:i:1 X1:i:
0 XM:i:0 X0:i:0 XG:i:0 MD:Z:38
ILLUMINA-GA_0000:1:1:4010:1065#0 147 scaffold_76 866076 60 38M = 865926 -188 ACTGTTATGTAATAGGACCATCTAATGTCACAAACA ehfhhhhhhhhhhhghghfhghhhghhhghhhhhhh XT:A:U NM:i:2 SM:i:37 AM:i:37 X0:i:1 X1:i:
0 XM:i:2 X0:i:0 XG:i:0 MD:Z:2C24A10
ILLUMINA-GA_0000:1:1:4093:1065#0 99 scaffold_57 479190 12 38M = 479354 202 AAATAACTAAGAAATTTGTACAAATTTCTAAATCTT affffgegggaaffccfd_ffcdfdffgfcgggfgggg XT:A:R NM:i:1 SM:i:0 AM:i:0 X0:i:7 X1:i:
2 XM:i:1 X0:i:0 XG:i:0 MD:Z:12G25
ILLUMINA-GA_0000:1:1:4093:1065#0 147 scaffold_57 479354 20 38M = 479190 -202 GTCTTGCATCTCTTGGCAACTTGTGTCTCTTGATT afdffffffbBfQf_gghchhfhfgghhhdhffiffd XT:A:U NM:i:0 SM:i:20 AM:i:0 X0:i:1 X1:i:
2 XM:i:0 X0:i:0 XG:i:0 MD:Z:38
ILLUMINA-GA_0000:1:1:6805:1068#0 99 scaffold_11 3541452 0 38M = 3541616 202 AATGCCATTATCTCTAAGTGTTTGTTCATCCAAAGTG hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh XT:A:R NM:i:1 SM:i:0 AM:i:0 X0:i:43 XM:i:
1 X0:i:0 XG:i:0 MD:Z:8A29

```

# Most common applications of NGS

## RNA-seq /Transcriptomics

- Quantitative
- Descriptive
  - Alternative splicing
- miRNA profiling

## ChIP-seq /Epigenomics

- Protein-DNA interactions
- Active transcription factor binding sites
- Histone methylation

## Resequencing

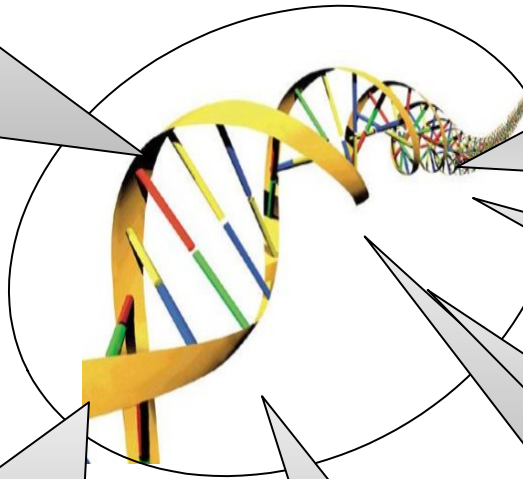
- Mutation calling
- Profiling
- Genome annotation

## *De novo* sequencing

## Exome sequencing Targeted sequencing

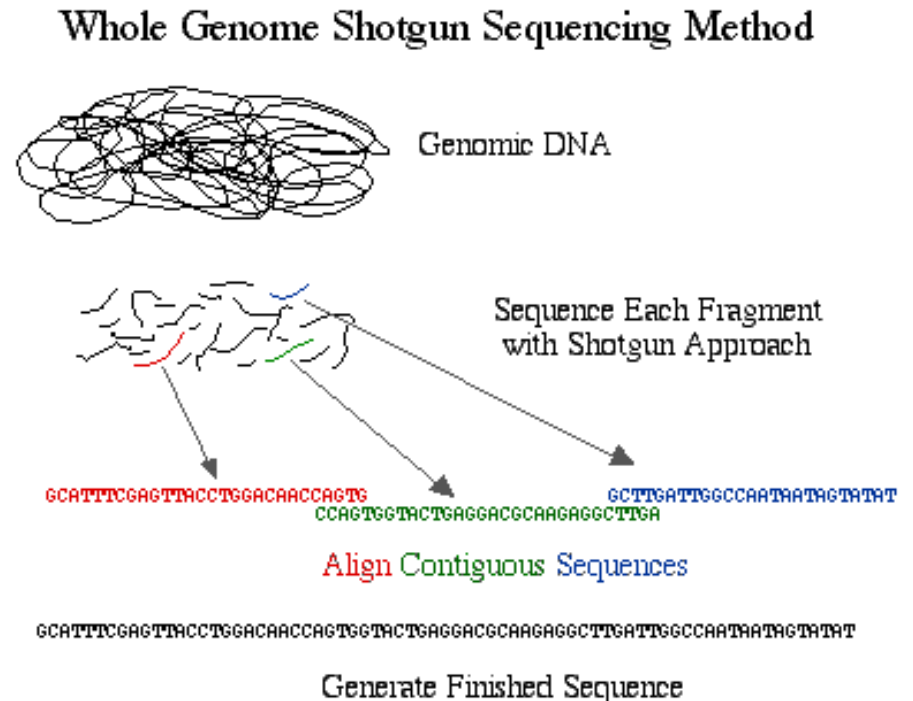
## Copy number variation

## Metagenomics Metatranscriptomics



# DNA sequencing - 1

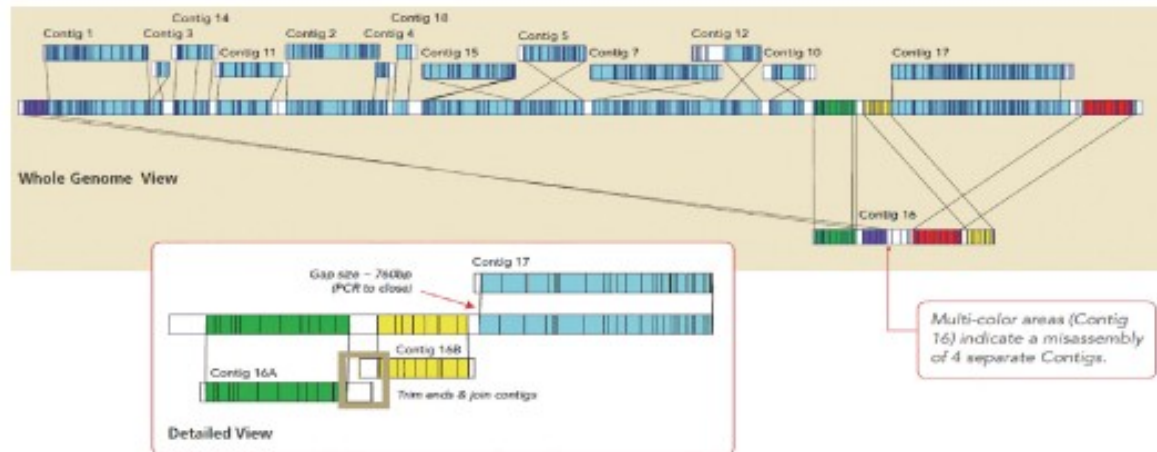
- **Whole GENOME Resequencing**
  - Need reference genome
  - Variation discovery



# DNA sequencing - 2

- **Whole GENOME “de novo” sequencing**

- Uncharacterized genomes with no reference genome available
- known genomes where significant structural variation is expected.
- Long reads or mate-pair libraries. Sequencing mostly done by Roche 454 and also Illumina.
- Assembly of reads is needed: Computational intensive
- E.g. Genome bacteria sequencing

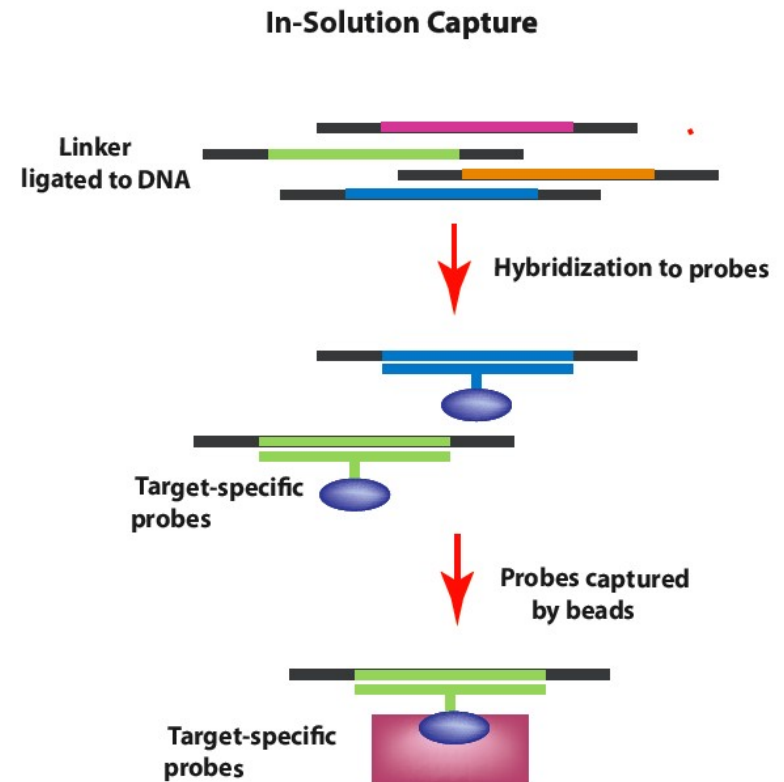


# DNA sequencing - 3

- **Whole EXOME Resequencing**

- Need reference genome
  - Available for Human and Mouse
- Variation discovery on ORFs
  - 2% of human genome (lower cost)
  - 85% disease mutation are in the exome
- Need probes complementary to exons
  - Nimblegen
  - Agilent

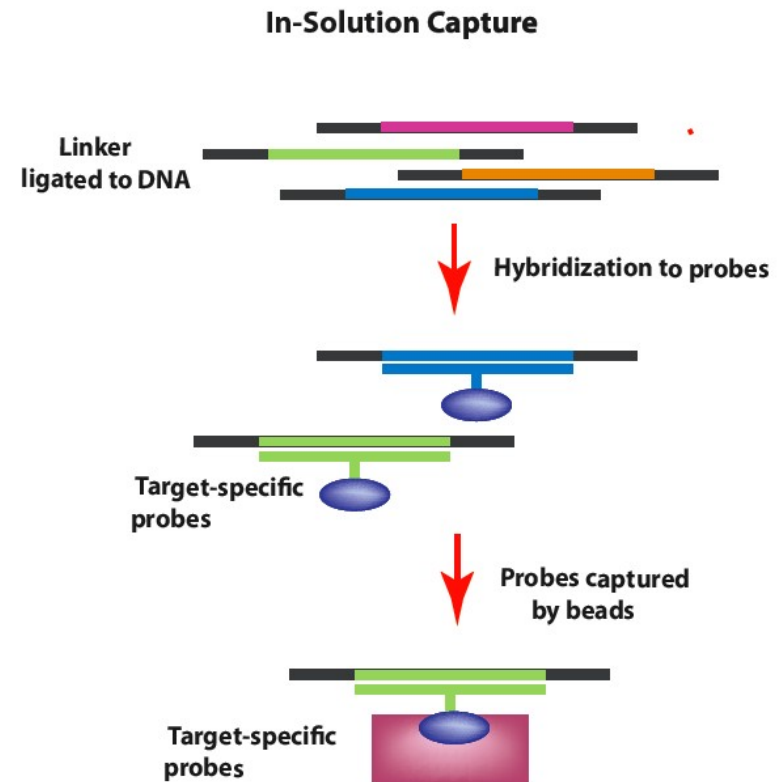
- E.g. Human exome





# DNA sequencing - 4

- **Targeted Resequencing**
  - Capture of specific regions in the genome
- **Custom genes panel sequencing**
  - Allows to cover high number of genes related to a disease
  - *E.g. Disease gene panel*
- Low cost and quicker than capillary sequencing
- Multiplexing is possible
- Need custom probes complementary to the genomic regions
  - Nimblegen
  - Agilent

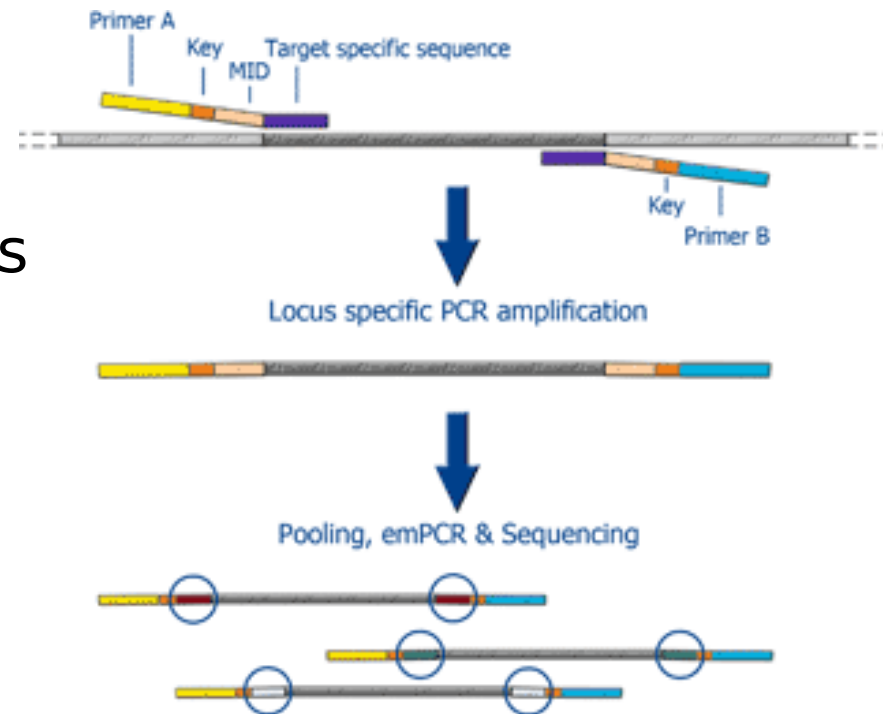


# DNA sequencing - 5

- **Amplicon sequencing**

- Sequencing of regions amplified by PCR.
- Shorter regions to cover than targeted capture
- No need of custom probes
- Primer design is needed
- High fidelity polymerase
- Multiplexing is needed

- *E.g. P53 exon amplicon sequencing*



# Transcriptomics - 1

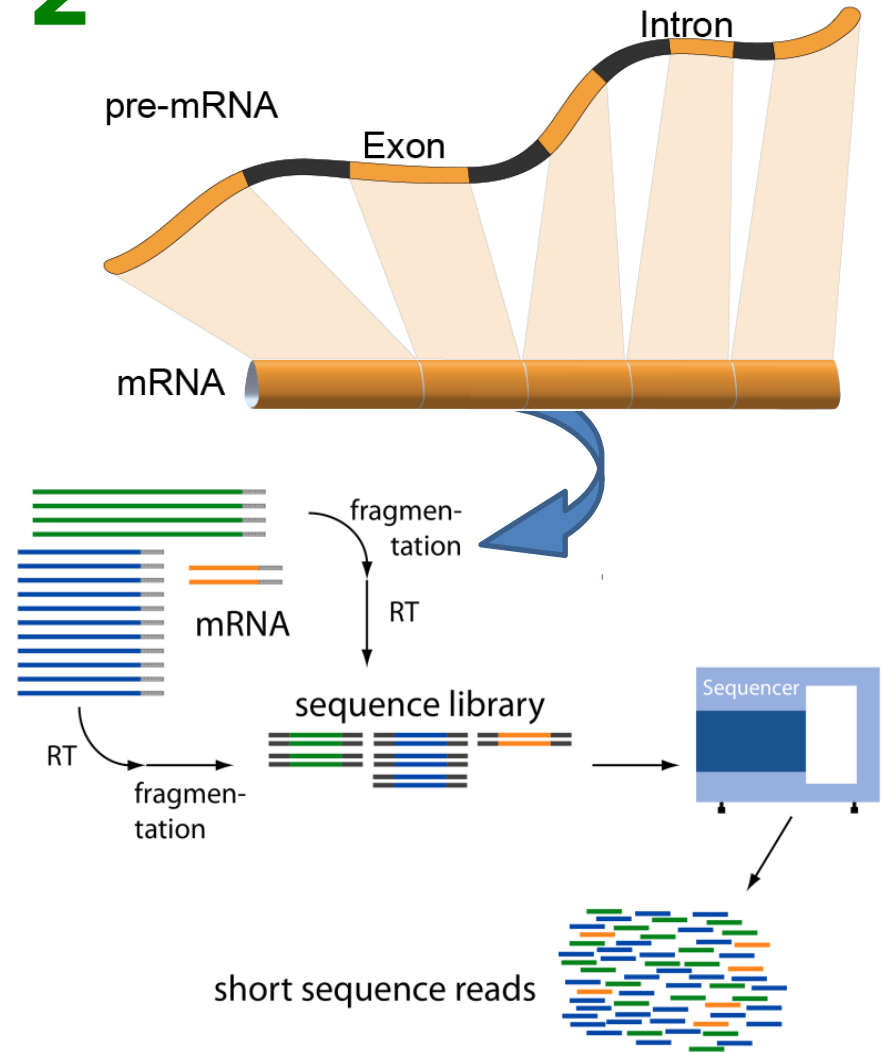
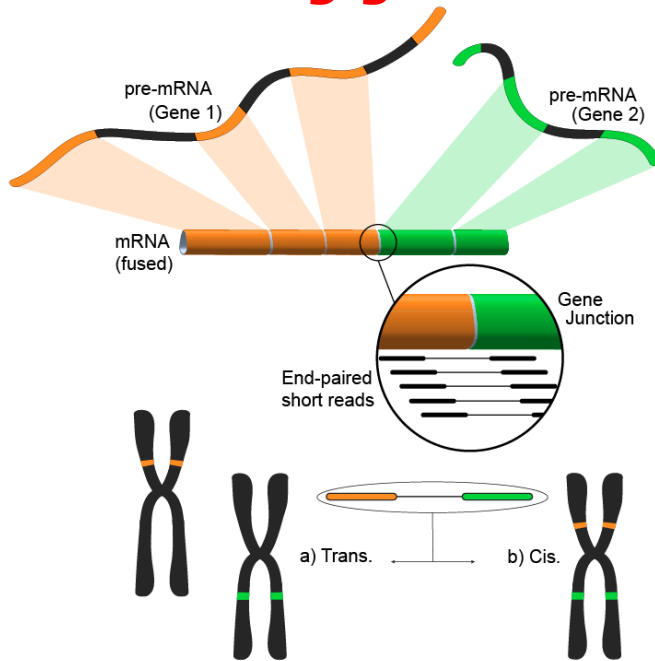
- **RNA-Seq**

- Sequencing of mRNA
- rRNA depleted samples
- Very high dynamic range
- No prior knowledge of expressed genes
- Gives information about (richer than microarrays)
  - Differential expression of **known or unknown** transcripts during a treatment or condition
  - **Isoforms** and
  - New **alternative splicing** events
  - **Non-coding** RNAs
  - Post-transcriptional mutations or **editing**,
  - **Gene fusions**.

# Transcriptomics - 2

- **RNA-Seq**

- Sequencing of **mRNA**
- **Detecting gene fusions**



# Applications of RNAseq

## Qualitative:

- \* Alternative splicing
- \* Antisense expression
- \* Extragenic expression
- \* Alternative 5' and 3' usage
- \* Detection of fusion transcripts

....

Tophat/Cufflinks  
Scripture  
Alexa

## Quantitative:

- \* Differential expression
- \* Dynamic range of gene expression

....

edgeR  
DESeq  
baySeq  
**NOISeq**

# Advantages of RNAseq?

## RNAseq

- \* Non targeted transcript detection
- \* No need of reference genome
- \* Strand specificity
- \* Find novel splicing sites
- \* Larger dynamic range
- \* Detects expression and SNVs
- \* Detects rare transcripts

....

## microarrays

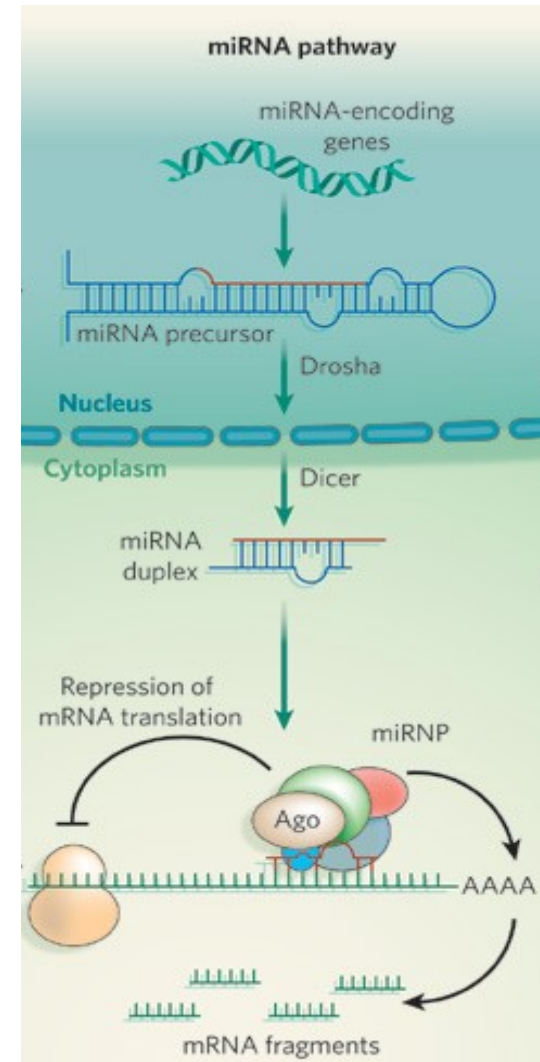
- \* Restricted to probes on array
- \* Needs genome knowledge
- \* Normally, not strand specific
- \* Exon arrays difficult to use
- \* Smaller dynamic range
- \* Does not provide sequence info
- \* Rare transcripts difficult

....

and.... are there any disadvantages?????

# Transcriptomics - 3

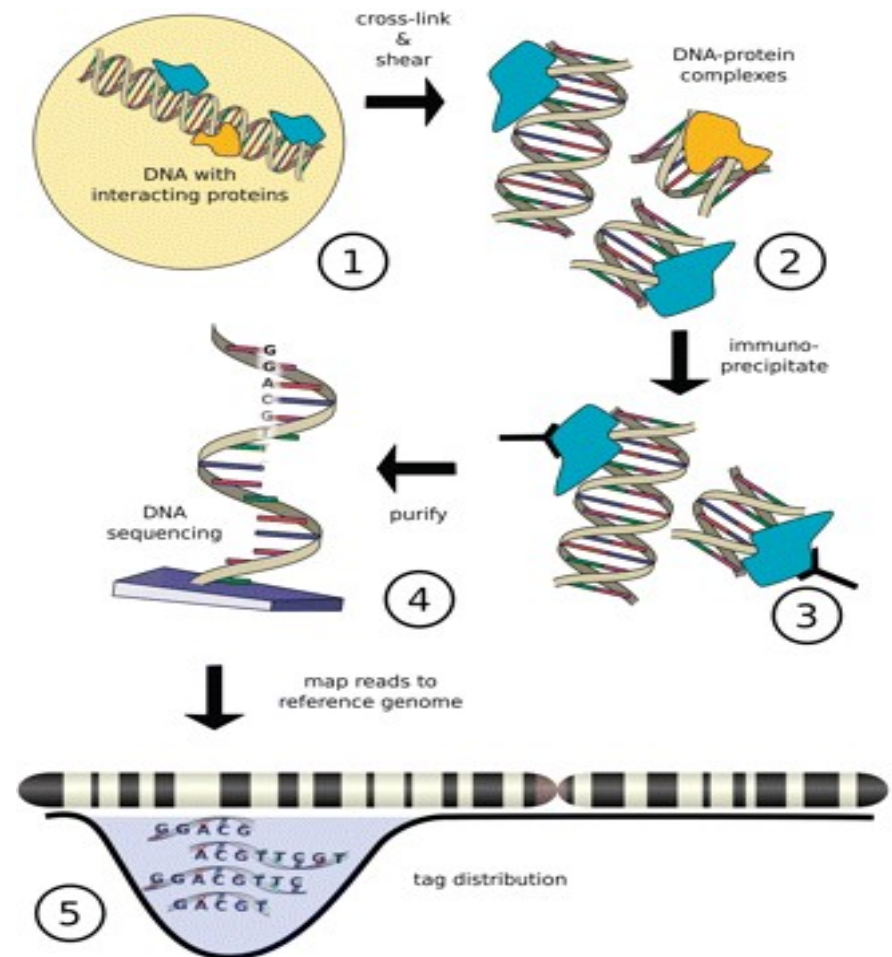
- **miRNA/small nonCoding RNA sequencing**
  - RNA Size selection step
    - 18-40 bp
  - Profiling of known miRNAs
  - miRNA discovery



# TFBS detection

## ChIP-Seq

- Identification of genomic region for gDNA binding proteins:
- Transcription Factor binding site detection





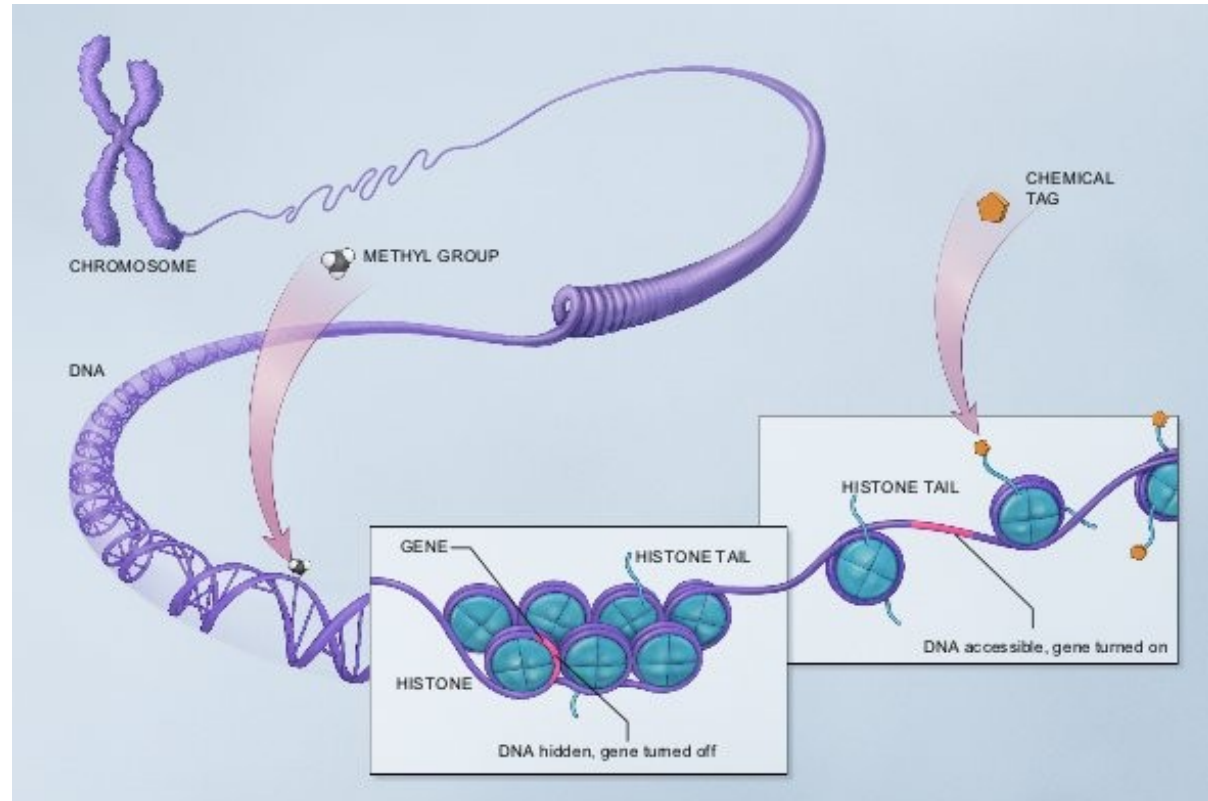
# Epigenomics - I

**Epigenomics** refers to functionally relevant modifications to the genome that do not involve a change in the nucleotide sequence

- *Play a role in turning genes off or on*

## Epigenomic Marks.

- Methyl groups attach to the backbone of a DNA molecule.
- A variety of chemical tags attach to the tails of histones. This action affects how tightly DNA is wound around the histones.



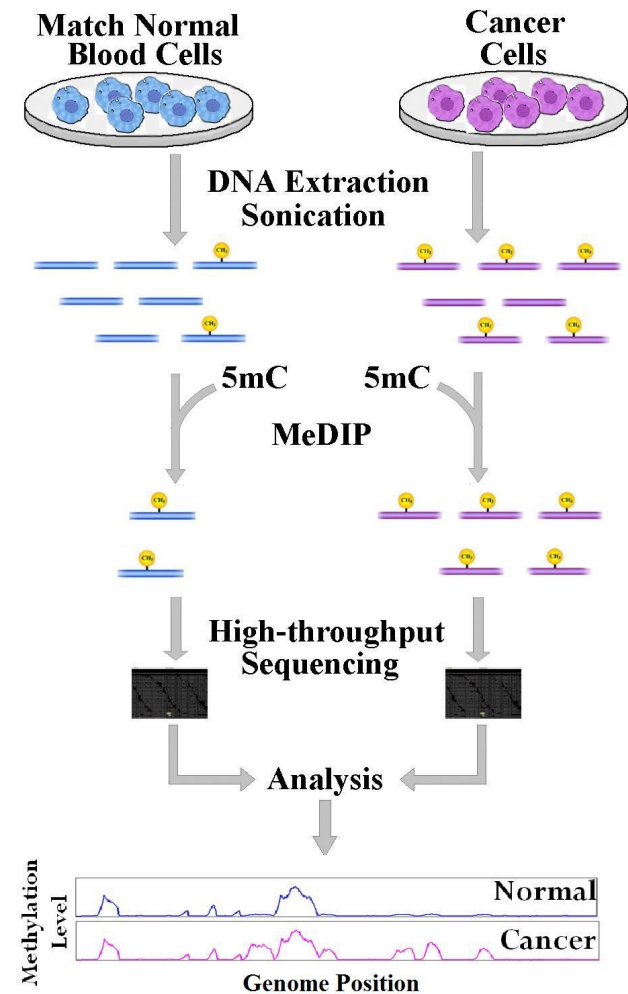
**ChIP-Seq:** Histone methylation detection

# Epigenomics - 2

- **Methyl-Seq**

- CpG island methylation
- Bisulfite sequencing-based method

- > E.g. Cancer studies.
  - Different degree of chromatin methylation affects expression of genes



**New huge projects coming**

- **Many big projects during the last years:**

- ENCODE <http://genome.ucsc.edu/ENCODE/>
- 1000 Genomes projects <http://www.1000genomes.org/>
- ICGC <http://icgc.org/>
- ...

- **New projects coming soon:**

- BRIDGE <https://bridgestudy.medschl.cam.ac.uk/index.shtml>
- Genomics England (<http://www.genomicsengland.co.uk/>) will produce tens of PB of data (1PB == 1000TB)

# **Successful NGStories**

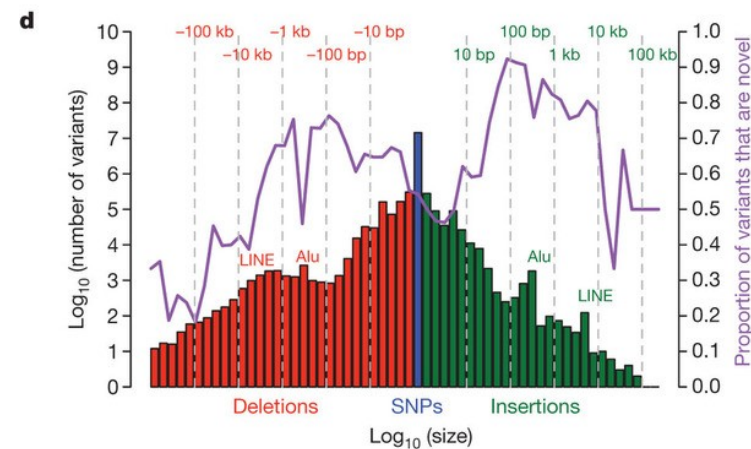
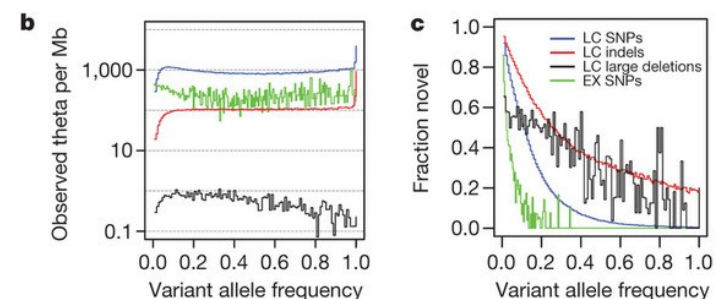
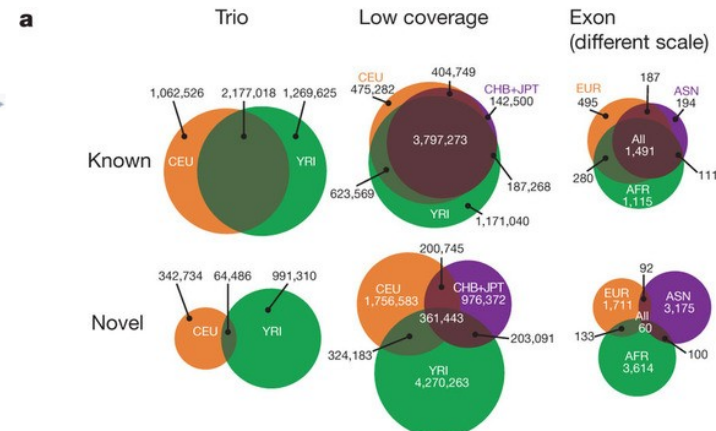
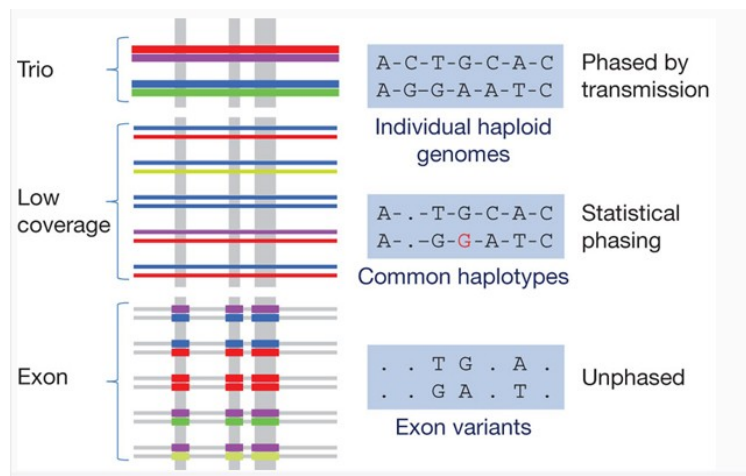
# A map of human genome variation from population-scale sequencing

The 1000 Genomes Project Consortium

Affiliations | Contributions | Corresponding author

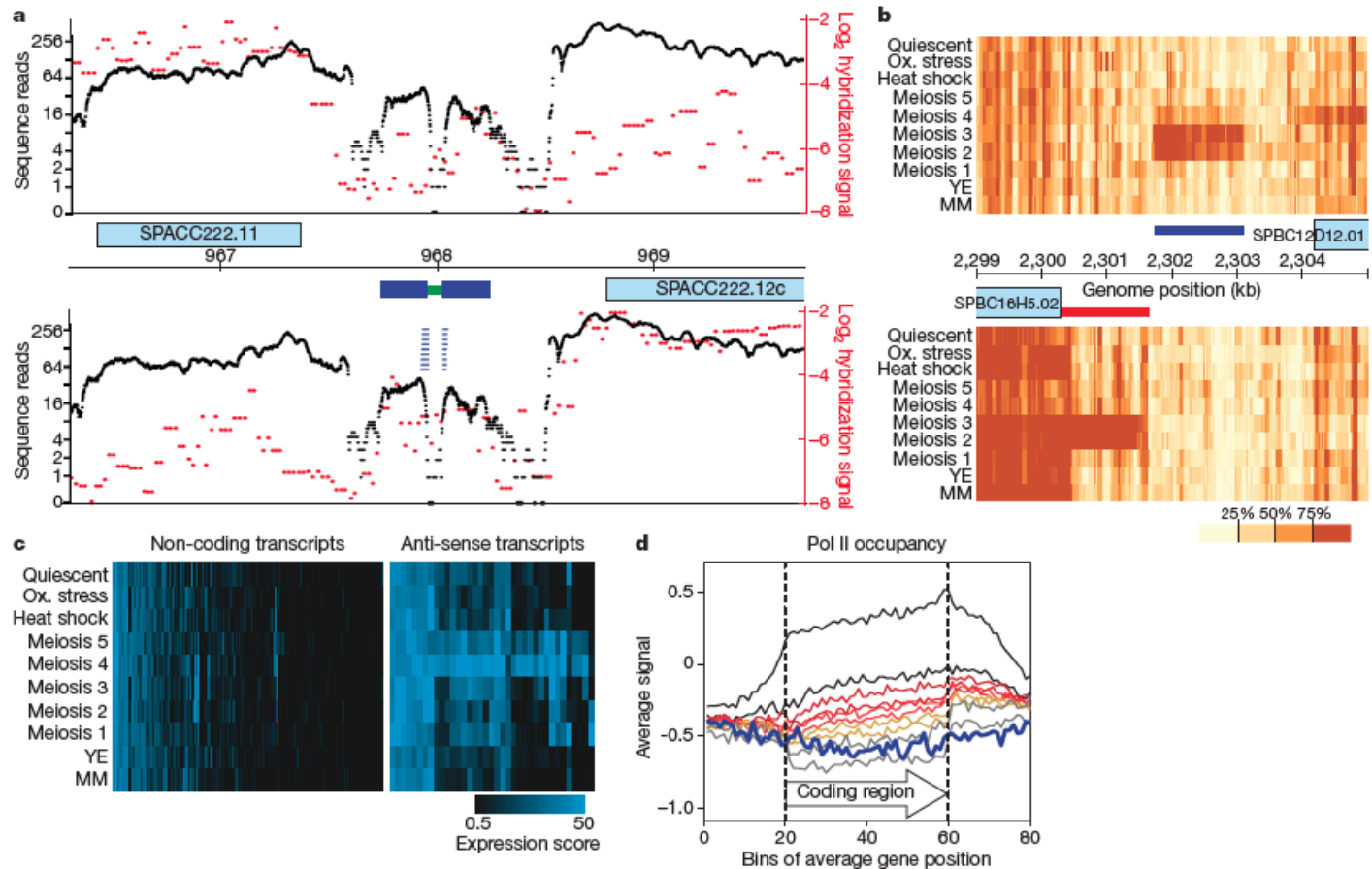
Nature 467, 1061–1073 (28 October 2010) | doi:10.1038/nature09534

Received 20 July 2010 | Accepted 30 September 2010 | Published online 27 October 2010



# Dynamic repertoire of a eukaryotic transcriptome surveyed at single-nucleotide resolution

Brian T. Wilhelm<sup>1\*†</sup>, Samuel Marguerat<sup>1\*†</sup>, Stephen Watt<sup>1†</sup>, Falk Schubert<sup>1†</sup>, Valerie Wood<sup>1</sup>, Ian Goodhead<sup>1†</sup>, Christopher J. Penkett<sup>1†</sup>, Jane Rogers<sup>1</sup> & Jürg Bähler<sup>1†</sup>



## Exome sequencing identifies the cause of a Mendelian disorder

Sarah B. Ng<sup>1,\*</sup>, Kati J. Buckingham<sup>2,\*</sup>, Choli Lee<sup>1</sup>, Abigail W. Bigham<sup>2</sup>, Holly K. Tabor<sup>2</sup>, Karin M. Dent<sup>3</sup>, Chad D. Huff<sup>4</sup>, Paul T. Shannon<sup>5</sup>, Ethylin Wang Jabs<sup>6,7</sup>, Deborah A. Nickerson<sup>1</sup>, Jay Shendure<sup>1,†</sup>, and Michael J. Bamshad<sup>1,2,8,†</sup>

<sup>1</sup>Department of Genome Sciences, University of Washington, Seattle, Washington, USA

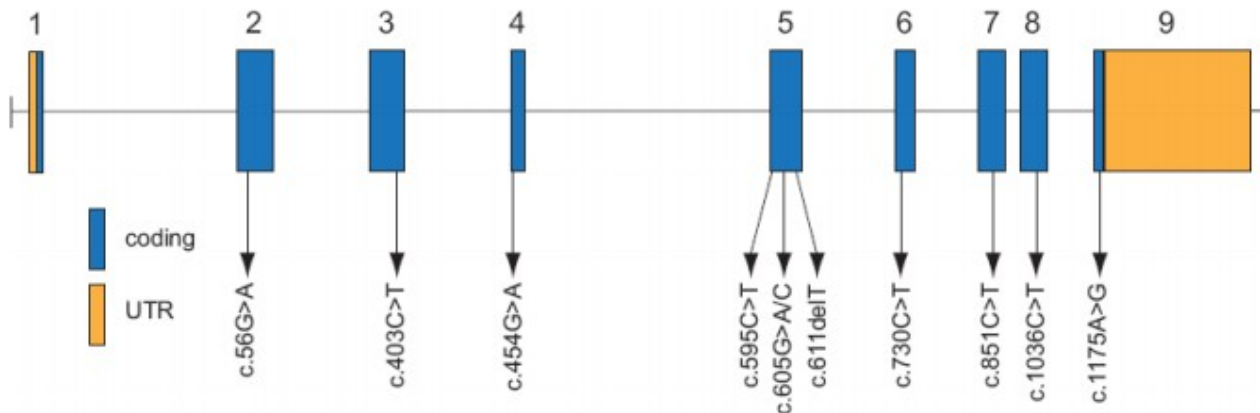
<sup>2</sup>Department of Pediatrics, University of Washington, Seattle, Washington, USA <sup>3</sup>Department of

Pediatrics, University of Utah, Salt Lake City, Utah, USA <sup>4</sup>Department of Human Genetics,

University of Utah, Salt Lake City, Utah, USA <sup>5</sup>Institute of Systems Biology, Seattle WA, USA

<sup>6</sup>Department of Genetics and Genomic Sciences, Mount Sinai School of Medicine, New York, New York, USA <sup>7</sup>Department of Pediatrics, Johns Hopkins University, Baltimore, Maryland <sup>8</sup>Seattle

Children's Hospital, Seattle, Washington, USA



**Figure 2. Genomic structure of the exons encoding the open reading frame of *DHODH***  
*DHODH* is composed of 9 exons that encode untranslated regions (orange) and protein coding sequence (blue). Arrows indicate the locations of 11 different mutations found in 6 families with Miller syndrome.



Miller syndrome



# Targeted next-generation sequencing of a cancer transcriptome enhances detection of sequence variants and novel fusion transcripts

Joshua Z Levin\*, Michael F Berger<sup>†</sup>, Xian Adiconis\*, Peter Rogov\*, Alexandre Melnikov\*, Timothy Fennell<sup>‡</sup>, Chad Nusbaum\*, Levi A Garraway<sup>†§</sup> and Andreas Gnirke\*

Addresses: \*Genome Sequencing and Analysis Program, Broad Institute of MIT and Harvard, 320 Charles Street, Cambridge, MA 02141, USA. <sup>†</sup>Cancer Program, Broad Institute of MIT and Harvard, 5 Cambridge Center, Cambridge, MA 02142, USA. \*Sequencing Platform, Broad Institute of MIT and Harvard, 320 Charles Street, Cambridge, MA 02141, USA. <sup>‡</sup>Department of Medical Oncology and Center for Cancer Genome Discovery, Dana-Farber Cancer Institute, Harvard Medical School, 44 Binney Street, Boston, MA 02115, USA.

<i>NUP214</i> (exon 29)	<i>XKR3</i> (exon 2)
caacctctgggttcagcttttgccaagcttcagCACCTGAGAATGGAGACAGTGTTTGAAGAGATGGATG	
T S G F S F C Q A S A P STOP	
<i>NUP214</i> (exon 29)	<i>XKR3</i> (exon 3)
caacctctgggttcagcttttgccaagcttcagGTGTTTGCACACCGTTAGAAATTACCACAAATGGTTGAAAAATC	
T S G F S F C Q A S G V C T P L E I T T N G STOP	
<i>NUP214</i> (exon 29)	<i>XKR3</i> (exon 4)
caacctctgggttcagcttttgccaagcttcagCATTGCTGATGACATTTCCCTGTTATCAGTTACTTATGGGGC	
T S G F S F C Q A S A L L M T F S L L S V T Y G	
<i>NUP214</i> (exon 27)	<i>XKR3</i> (exon 4)
atcttctccatcaggCATTGCTGATGACATTTCCCTGTTATCAGTTACTTATGGGGCCATTCGCTGCAATATACT	
F S P S G I A D D I F P V I S Y L W G H S L Q Y T	

**Figure 3**  
Sequences from *NUP214*-*XKR3* fusion transcripts detected after hybrid selection. After hybrid selection, 152 reads were aligned to the transcriptome and detected as *NUP214*-*XKR3* fusions. From top to bottom, we observed 137, four, eight, and three reads for these transcripts. The *NUP214* (exon 27) to *XKR3* (exon 4) has a stop codon downstream (not shown). Only *NUP214* (exon 29) to *XKR3* (exon 4) retains an open reading frame downstream of the fusion. Before hybrid selection, eight reads were aligned to the transcriptome and detected as *NUP214*-*XKR3* fusions; only the *NUP214* (exon 29) to *XKR3* (exon 2) transcript was detected. Sequence from *NUP214* DNA is shown as lower case, and from *XKR3*, as bold and upper case.

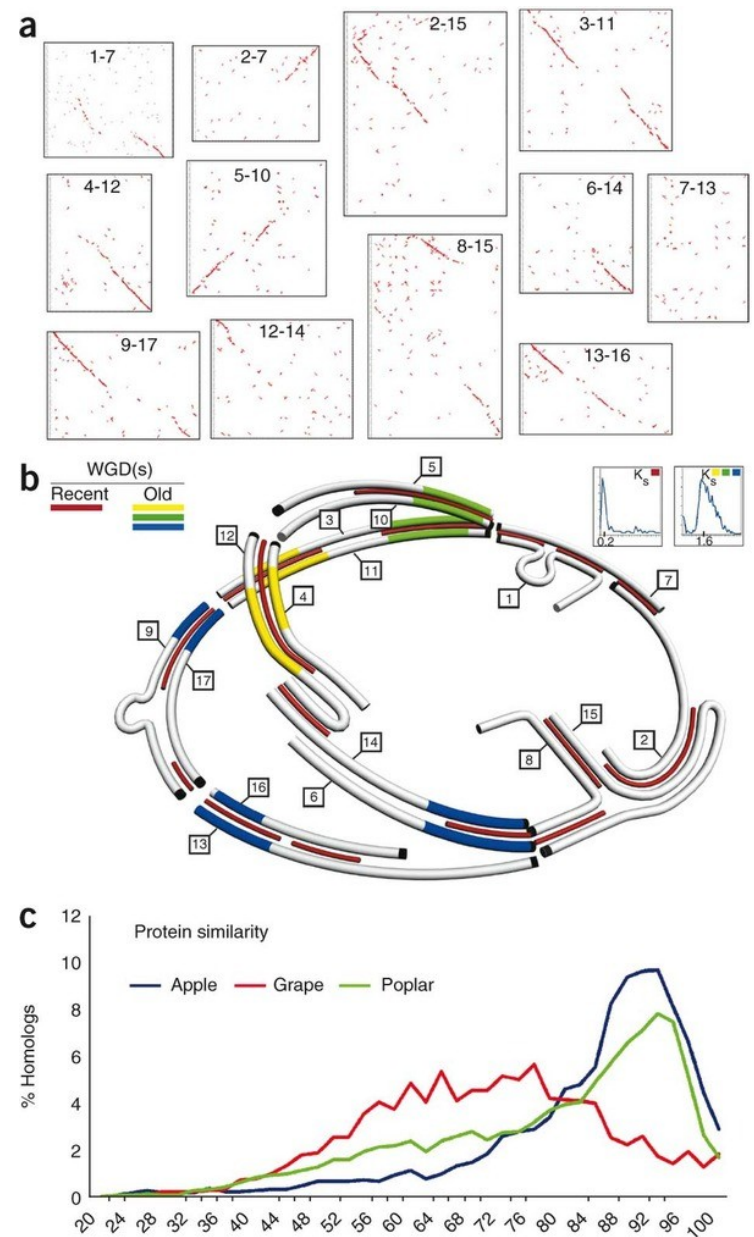
# The genome of the domesticated apple (*Malus × domestica* Borkh.)

Riccardo Velasco, Andrey Zharkikh, Jason Affourtit, Amit Dhingra, Alessandro Cestaro, Ananth Kalyanaraman, Paolo Fontana, Satish K Bhatnagar, Michela Troggio, Dmitry Pruss, Silvio Salvi, Massimo Pindo, Paolo Baldi, Sara Castelletti, Marina Cavauiuolo, Giuseppina Coppola, Fabrizio Costa, Valentina Cova, Antonio Dal Ri, Vadim Goremykin, Matteo Komjanc, Sara Longhi, Pierluigi Magnago, Giulia Malacarne, Mickael Malnoy *et al.*

Affiliations | Contributions | Corresponding author

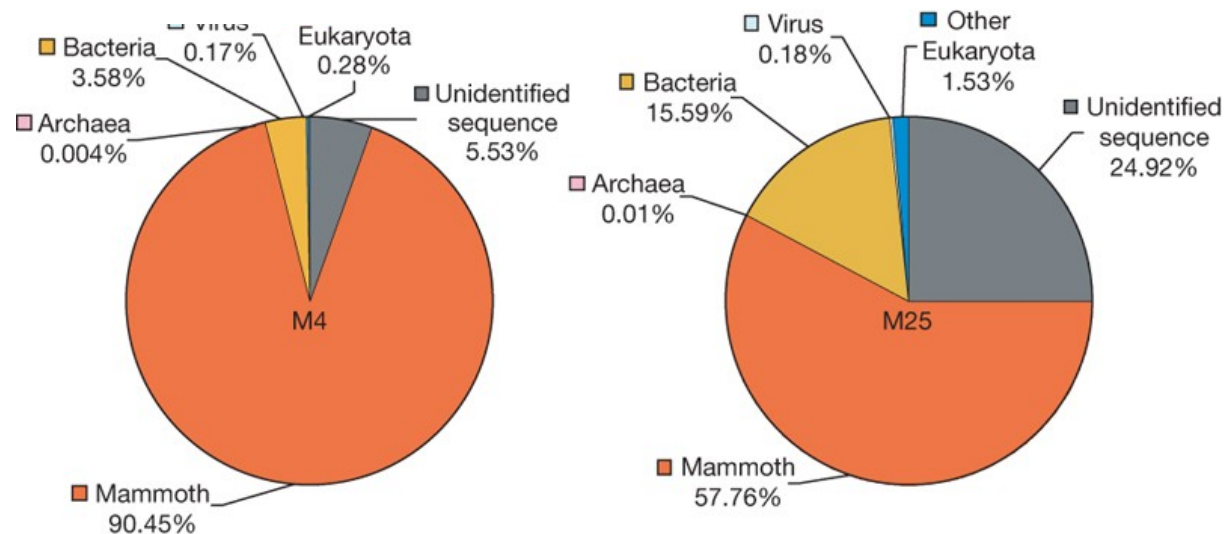
*Nature Genetics* **42**, 833–839 (2010) | doi:10.1038/ng.654

Received 19 November 2009 | Accepted 03 August 2010 | Published online 29 August 2010



## Sequencing the nuclear genome of the extinct woolly mammoth

Webb Miller<sup>1</sup>, Daniela I. Drautz<sup>1</sup>, Aakrosh Ratan<sup>1</sup>, Barbara Pusey<sup>1</sup>, Ji Qi<sup>1</sup>, Arthur M. Lesk<sup>1</sup>, Lynn P. Tomsho<sup>1</sup>, Michael D. Packard<sup>1</sup>, Fangqing Zhao<sup>1</sup>, Andrei Sher<sup>2,9</sup>, Alexei Tikhonov<sup>3</sup>, Brian Raney<sup>4</sup>, Nick Patterson<sup>5</sup>, Kerstin Lindblad-Toh<sup>5</sup>, Eric S. Lander<sup>5</sup>, James R. Knight<sup>6</sup>, Gerard P. Irzyk<sup>6</sup>, Karin M. Fredrikson<sup>7</sup>, Timothy T. Harkins<sup>7</sup>, Sharon Sheridan<sup>7</sup>, Tom Pringle<sup>8</sup> & Stephan C. Schuster<sup>1</sup>



Species composition of metagenomic DNA extracted from mammoth hair

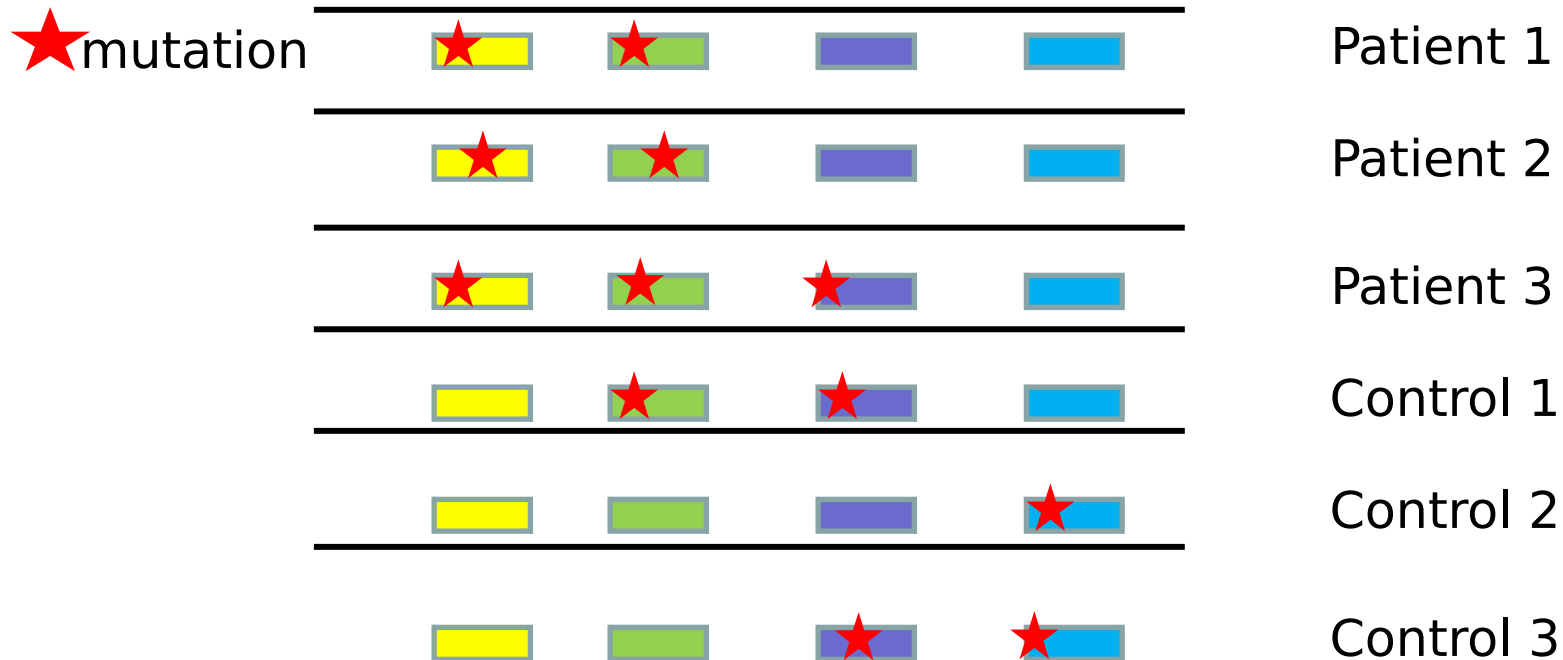
nature

**Not that easy, some challenges**

# The simplest case: monogenic disease due to a single gene



# The principle: comparison of patients (or families) and reference controls



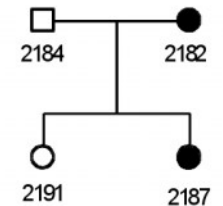
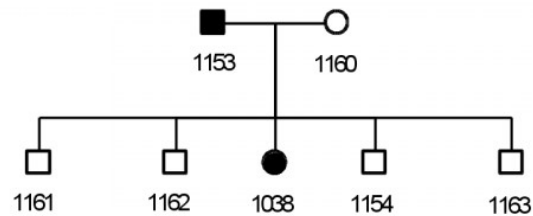
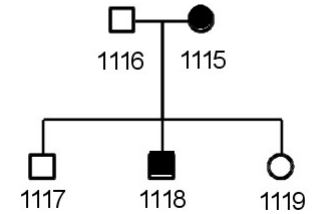
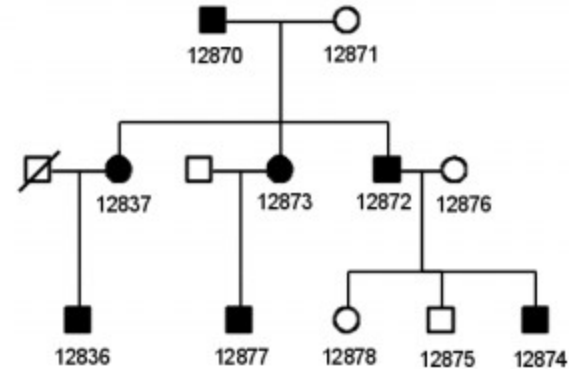
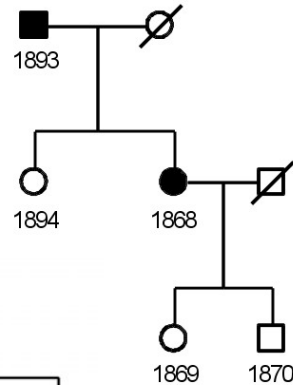
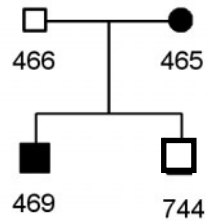
 **candidate gene** (shares mutation for all patients but no controls)

**Is this approach realistic?**

**Can we detect such rare variants so easily?**

- a) Interrogating 50Mb produces too many variants
- b) In many cases we are not hunting new but known variants
- c) Same phenotype can be due to different mutations and different genes

# Filtering with multiple family information



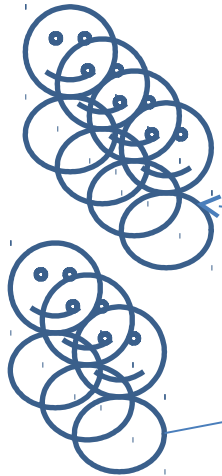
	Families					
	1	2	3	4	5	6
Variants	3403	82	4	0	0	0
Genes	2560	331	35	8	1	0

**Problem: how to prioritize putative candidate genes**



# Controls

# Cases



Many cases have to be used to obtain significant associations to many markers.

The only common element is the pathway (yet unknown) affected.

# Conclusions

NGS is  
revolutionizing how  
we do genome  
research

**But it will also  
revolutionize  
our lives....**

**If we manage  
to process and  
analyze ALL  
the DATA**

