

# Sequencing Technologies and Applications

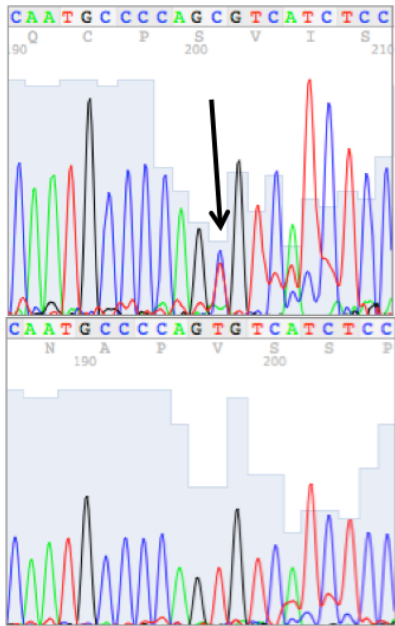
David K. Crossman, Ph.D.  
Department of Genetics  
Heflin Center for Genomic Sciences

# Genotyping Technologies

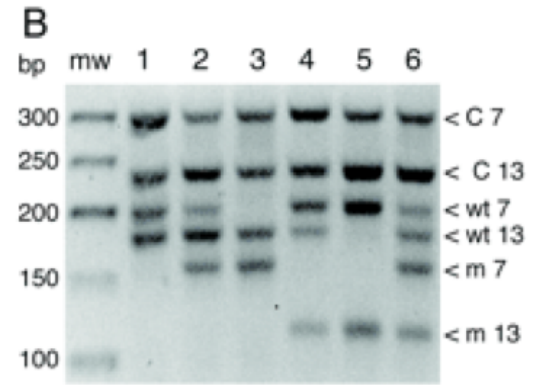
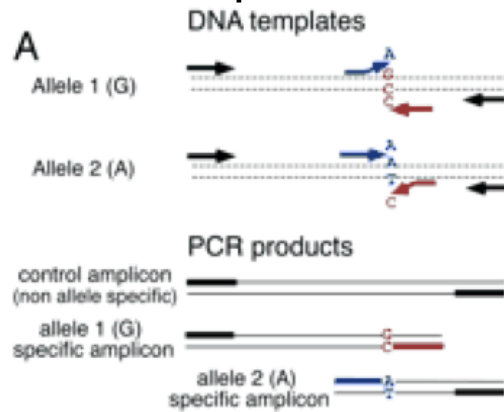
- The technology used depends on the number of variants and number of samples.
- Single to a few SNPs in a small population ( $\leq 960$ )
  - TaqMan
  - Pyrosequencing
  - Allele Specific PCR
- Intermediate # of SNPs, Intermediate population size
  - No good option here (BeadExpress or GoldenGate from Illumina but these are being discontinued)
- Large # of SNPs, large population
  - Infinium from Illumina (up to 5M SNPs per slide)

# Single SNP Analysis

## Direct Sequencing

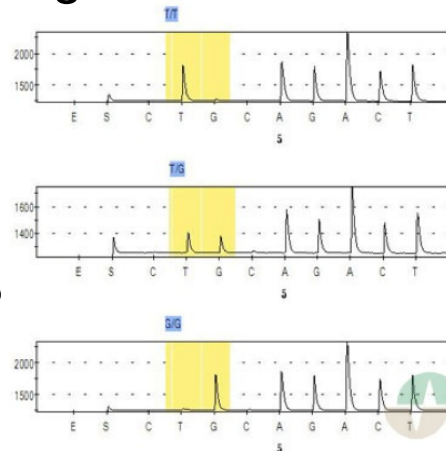
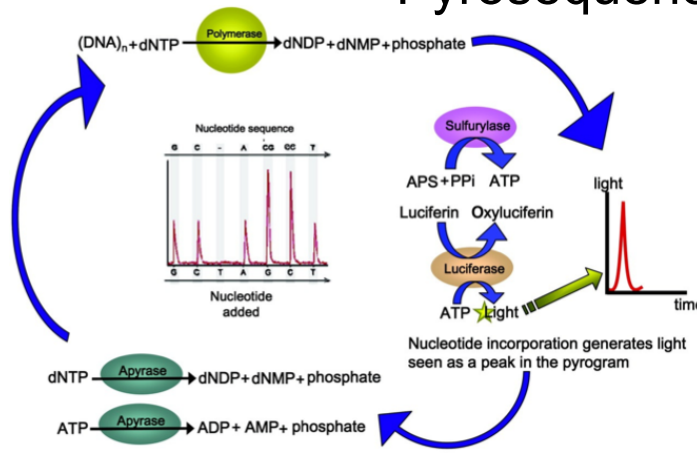


## Allele Specific PCR

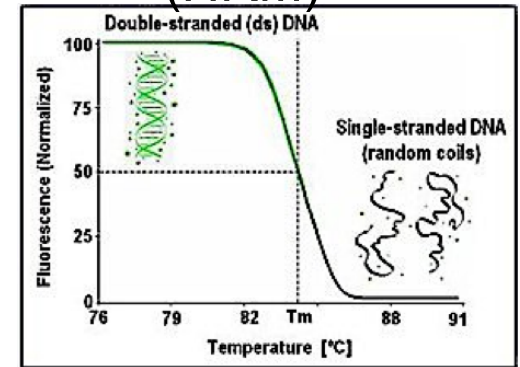


Piccioli, P. et al. Clinical Chem 2006 52:4:739

## Pyrosequencing



## High Resolution Melting (HRM)



[https://dna.utah.edu/Hi-Res/TOP\\_Hi-Res%20Melting.html](https://dna.utah.edu/Hi-Res/TOP_Hi-Res%20Melting.html)

# Cycle Sequencing with Dye Termination

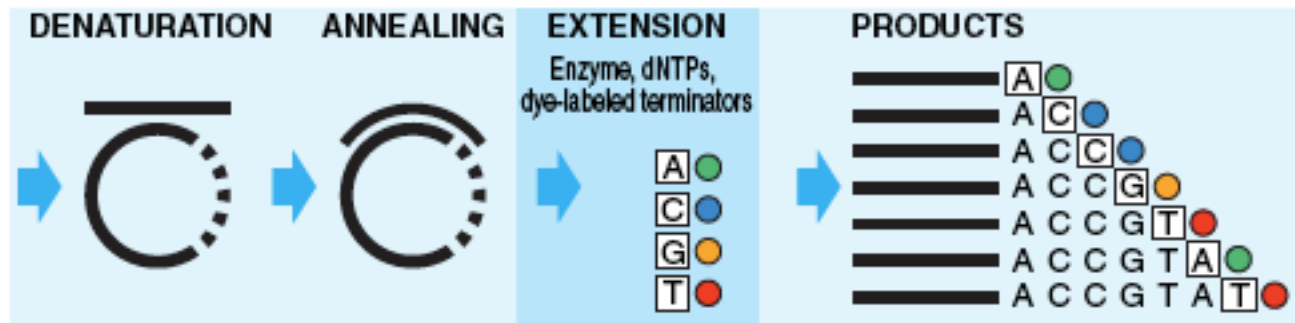
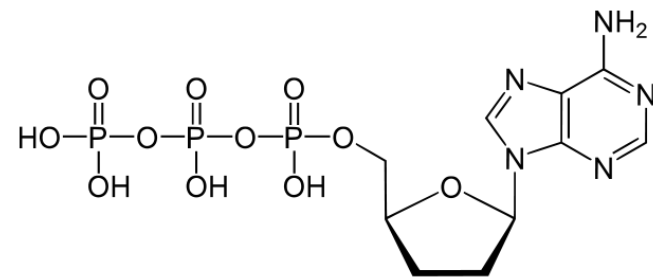
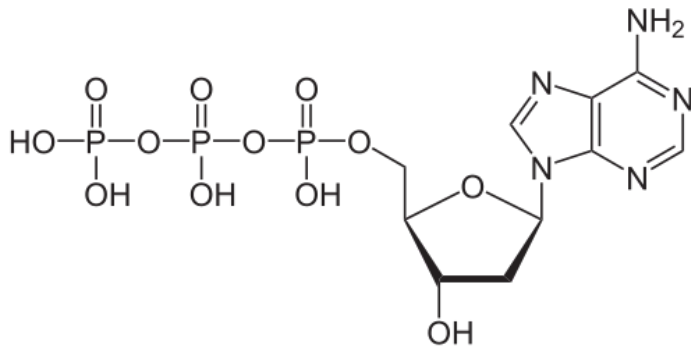
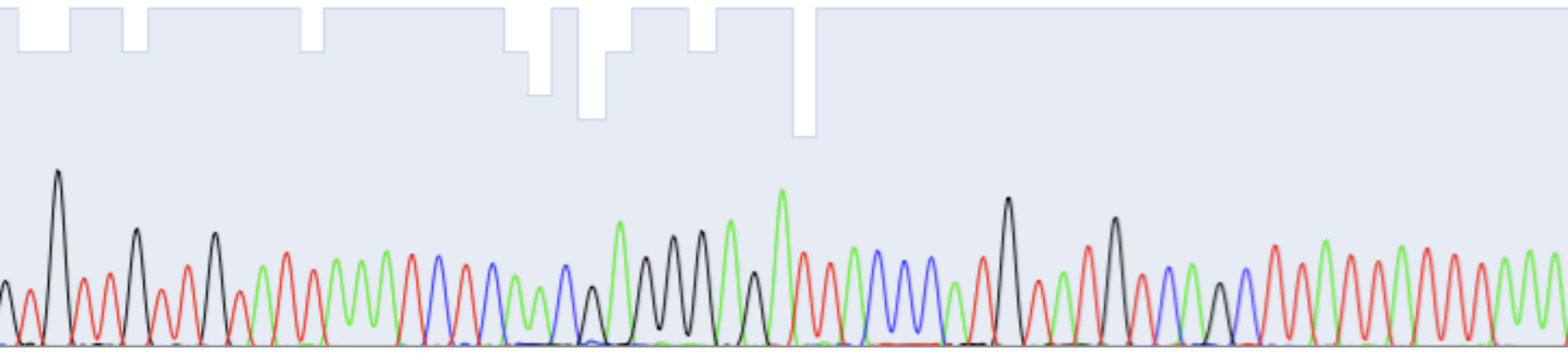


Figure 5 Diagram of dye terminator cycle sequencing



G T G T T G T T G T A T T A A A T C T C A A C G A G G G A G A T T A C C C A T G T A T G T C A G C T T A T T A T T T A A A  
C C C I K S 150 Q R G 160 R L P M 170 Y V S 180 L L F 190 K



# The Next Generation

# Illumina Platforms



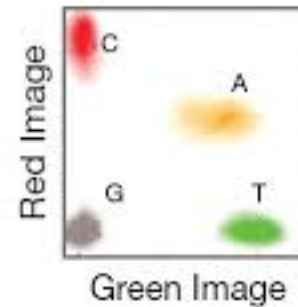
(

MiSeq  
 ~ Single flowcell  
 36bp increments  
 25-600bp per base seq  
 Single read and Paired end reads



HiSeq2500  
 Two flowcells  
 ~600billion bases sequenced  
 50bp-100bp increments  
 Lower cost per base sequenced  
 Single reads and Paired end reads  
 Rapid Runs 26-48hrs

# Illumina NextSeq 500



**ONE SYSTEM, TWO OUTPUT MODES**

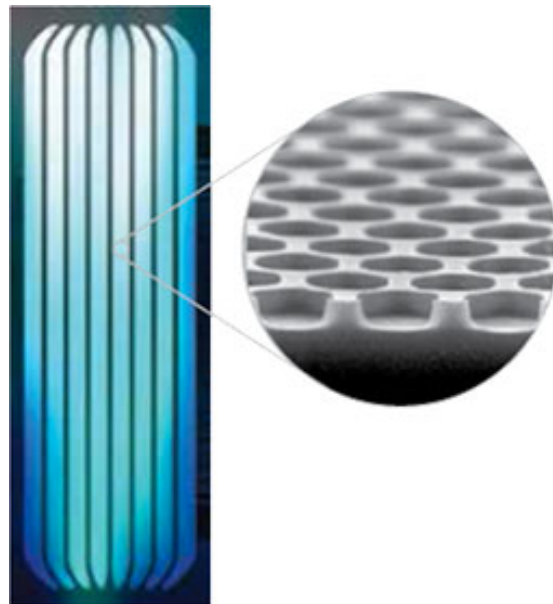
High-Output	Mid-Output
Up to 120 Gb 400M clusters FF 1 x 75 bp to 2 x 150 bp	Up to 40 Gb 130M clusters FF 2 x 75 bp to 2 x 150 bp
30x genome	2-3 assays
6-12 assays RNA-Seq	2-4 samples RNA-Seq
20 GEX profiles MSPT	6-36 panels

51 © 2018 Illumina Inc. All Rights Reserved.

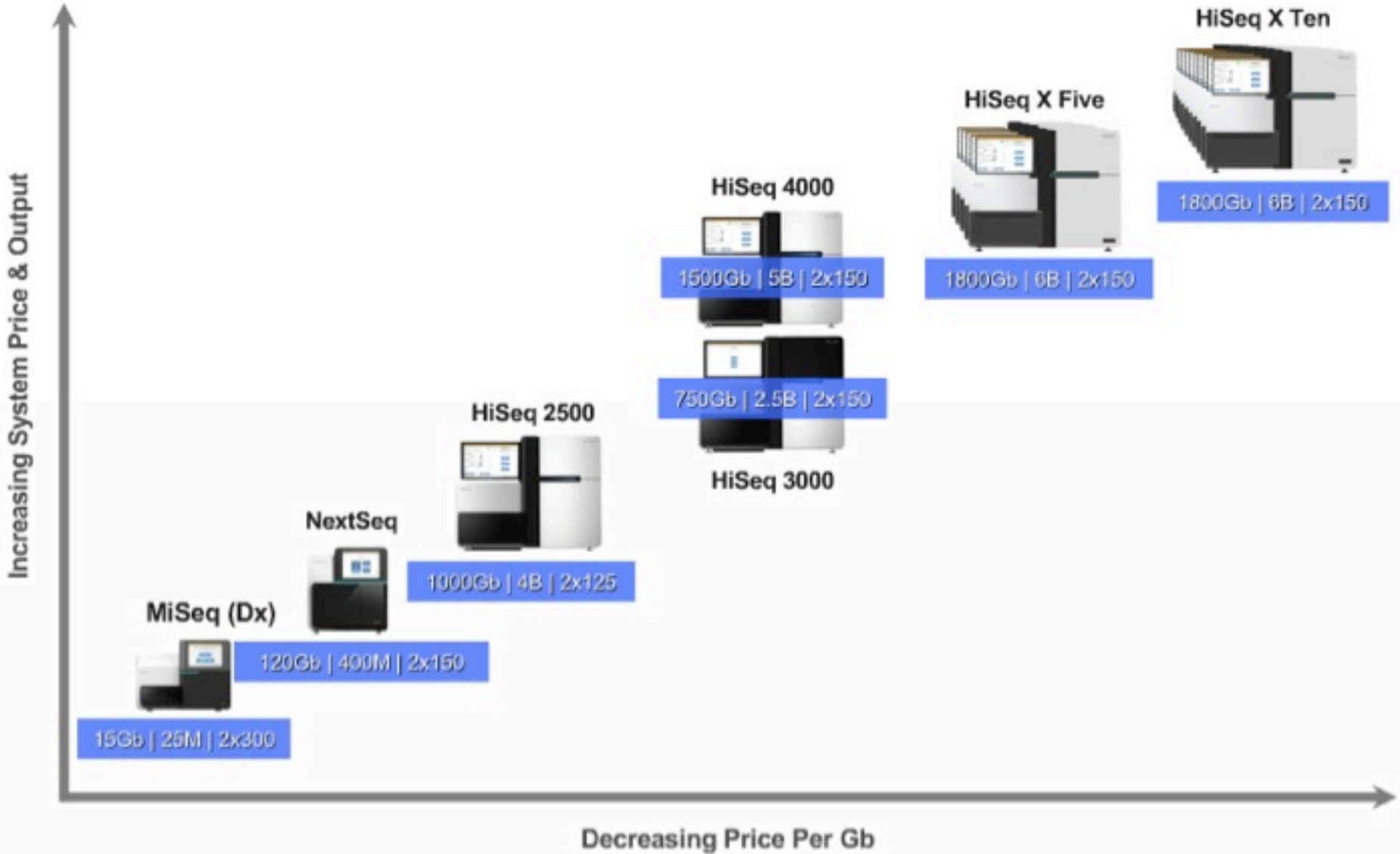
Accelerated detection of all four DNA bases is performed on the NextSeq 500 System using only two images to capture red and green filter wavelength bands. A bases will be present in both images (yellow cluster), C bases in red only, T bases in green only, and G bases in neither.



# Illumina HiSeq X Ten



# Sequencing Power For Every Scale.



# Useful Next-Gen Terms

- Cluster
  - Individual island of DNA molecules representing a single, unique template
- Clusters Passing filter
  - Number of clusters able to be distinguished by the software as individuals
- Fastq
  - DNA Sequence file that is able to be read by downstream analysis applications
- Q-Score
  - A quality score based on the Phred score from Sanger Sequencing which is the probability a base is incorrect at a give position. Example: Q30 means there is a 1:1000 chance the base is incorrect. Or stated another way it means the base call is 99.9% accurate
- Phasing/Prephasing
  - When the DNA sequencing reaction is either a base ahead or a base behind the majority of the other molecules
- Depth of Coverage
  - The average number of times a base is read within the genome
- Reads
  - Actual sequence

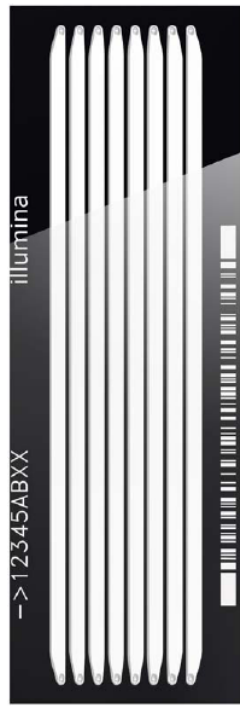
# Flowcells through time

2005



GAIIx

2010



HiSeq



2013 →



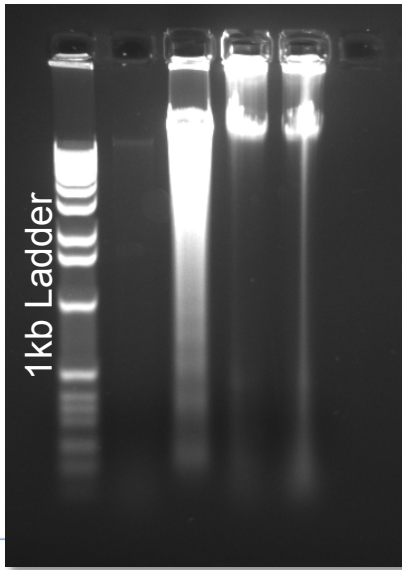
MiSeq



HiSeqX

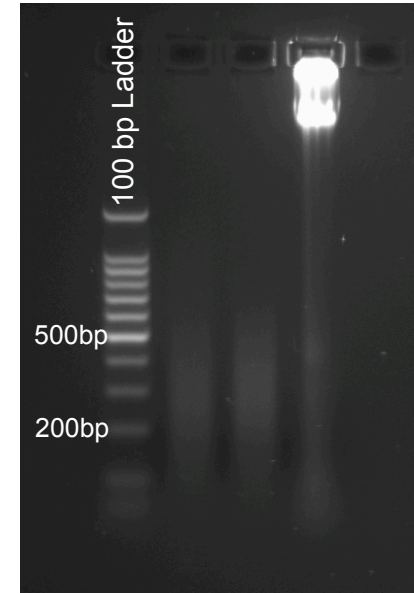
Not to scale

# DNA Library Prep and Flow cell Production

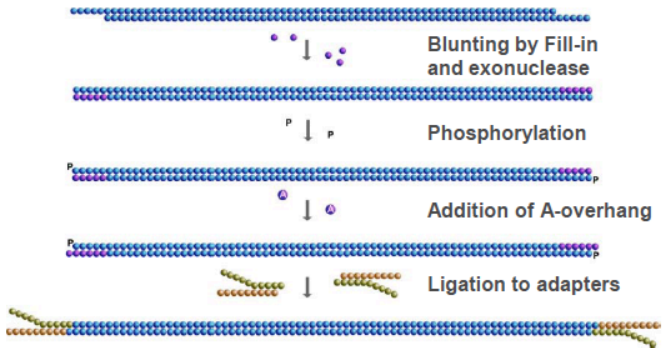


**S-series**

- Manual
- Single



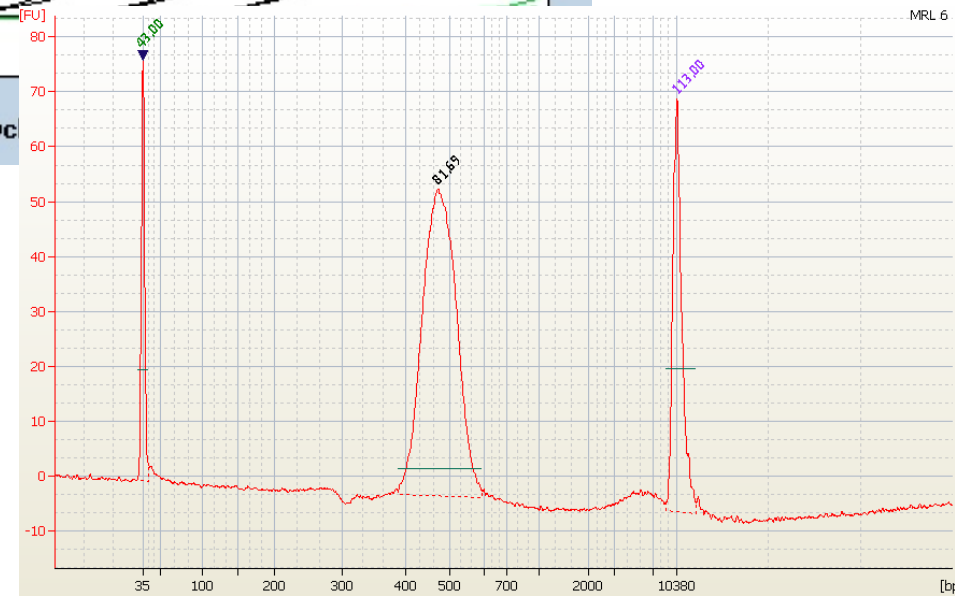
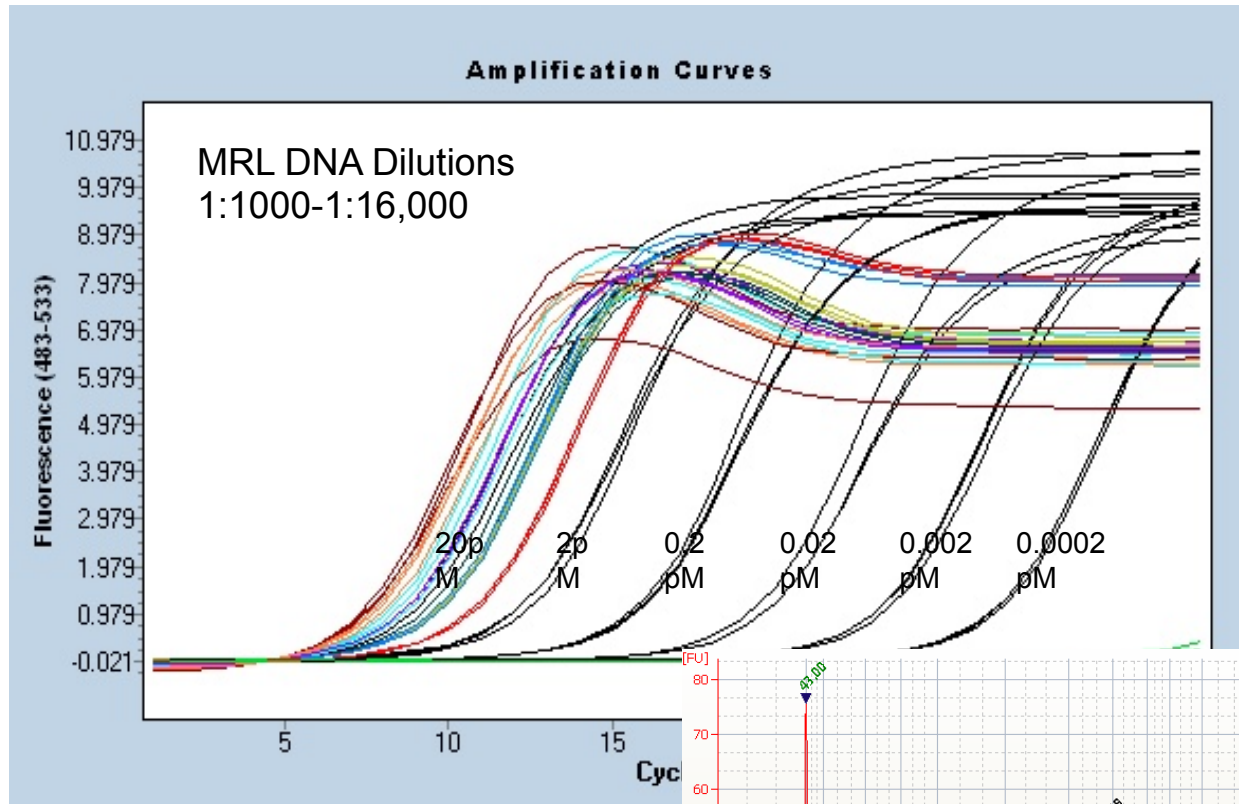
DNA fragments



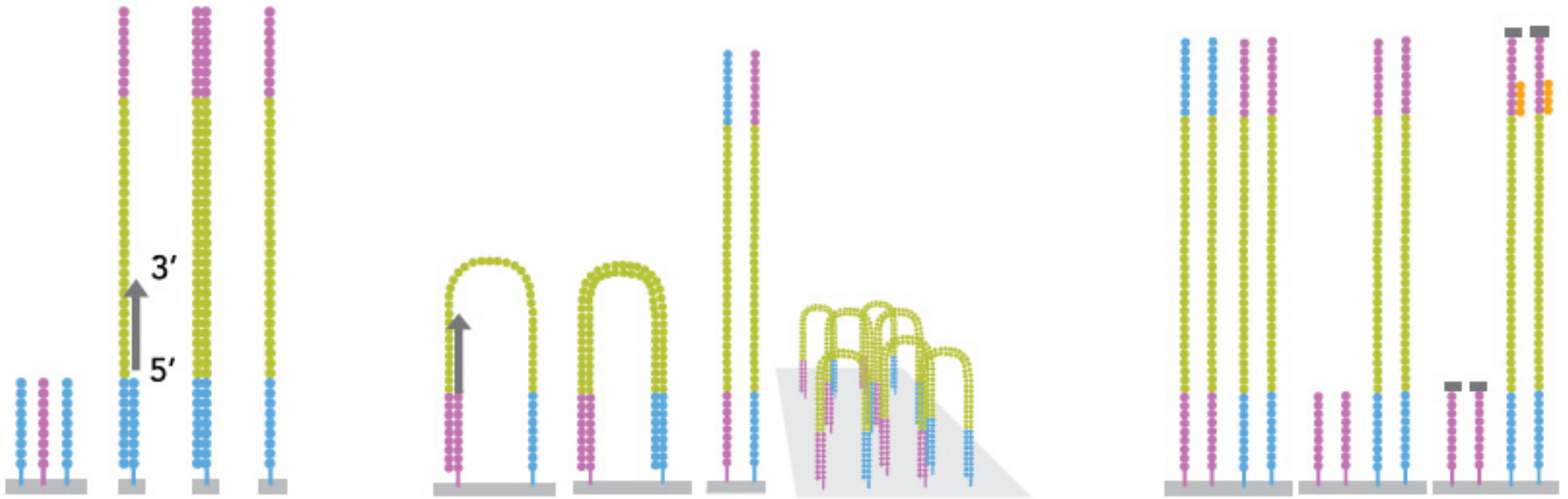
Version 3 HiSeq Flow Cell

Illumina's Library Preparation Workflow

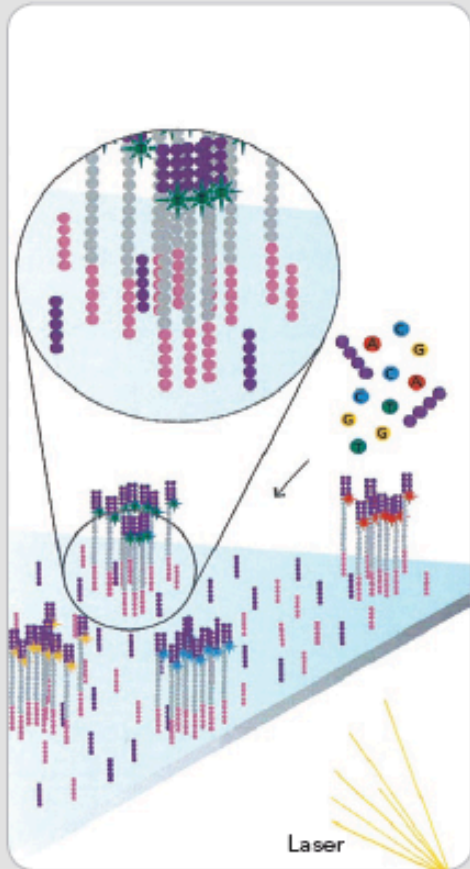
# Library Assessment and Quantitation



# Illumina Cluster Generation

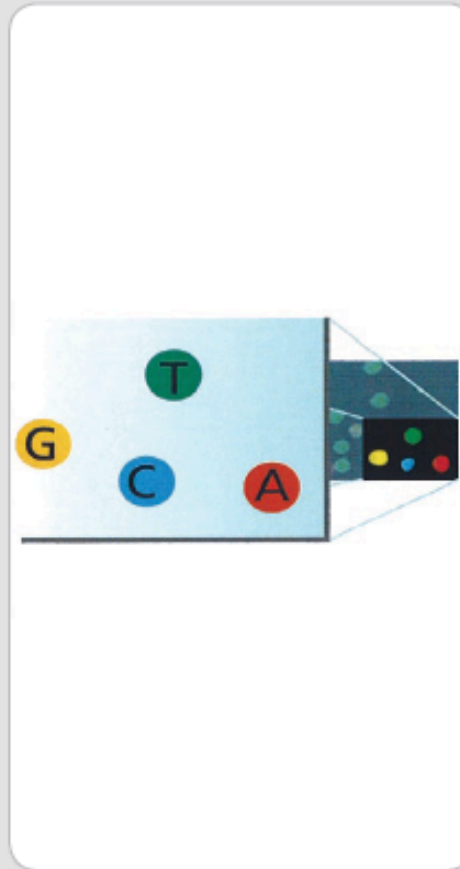


## 7. DETERMINE FIRST BASE



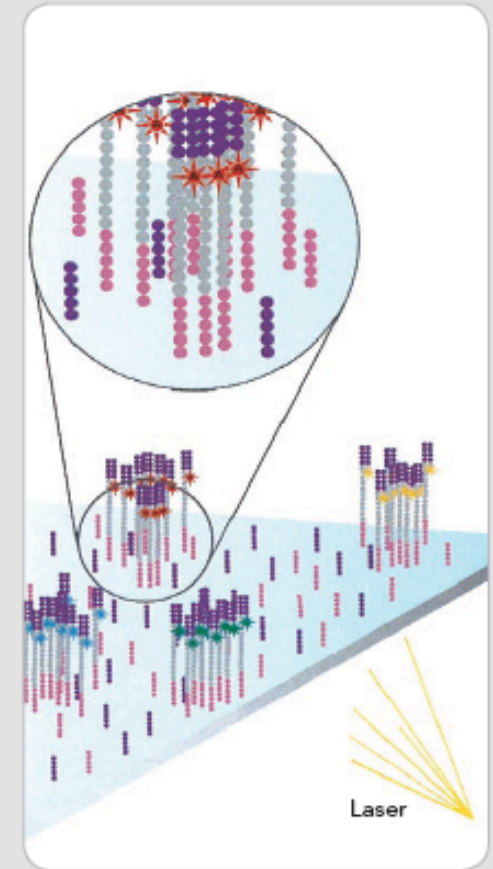
The first sequencing cycle begins by adding four labeled reversible terminators, primers, and DNA polymerase.

## 8. IMAGE FIRST BASE



After laser excitation, the emitted fluorescence from each cluster is captured and the first base is identified.

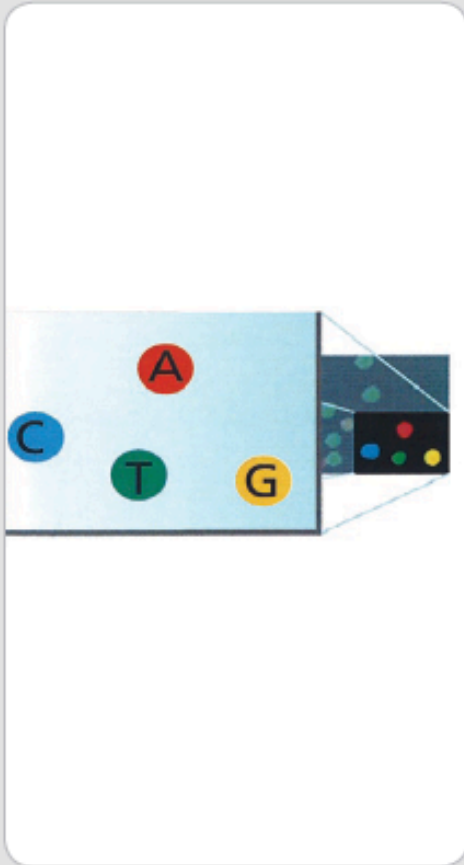
## 9. DETERMINE SECOND BASE



The next cycle repeats the incorporation of four labeled reversible terminators, primers, and DNA polymerase.

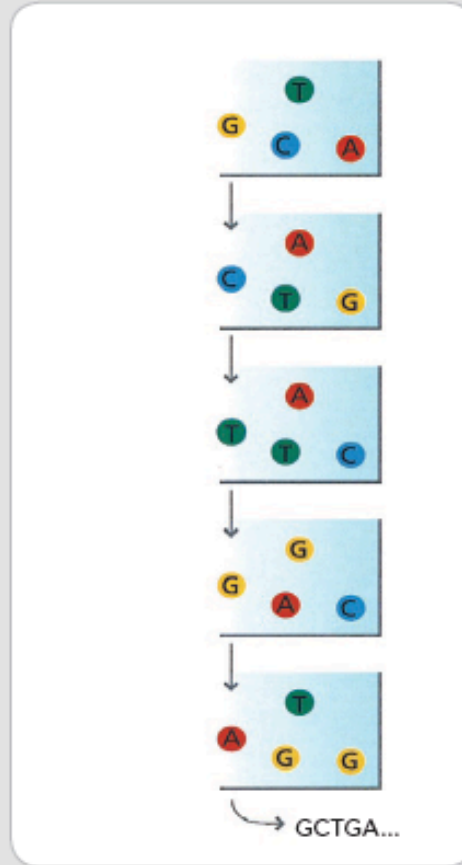


## 10. IMAGE SECOND CHEMISTRY CYCLE



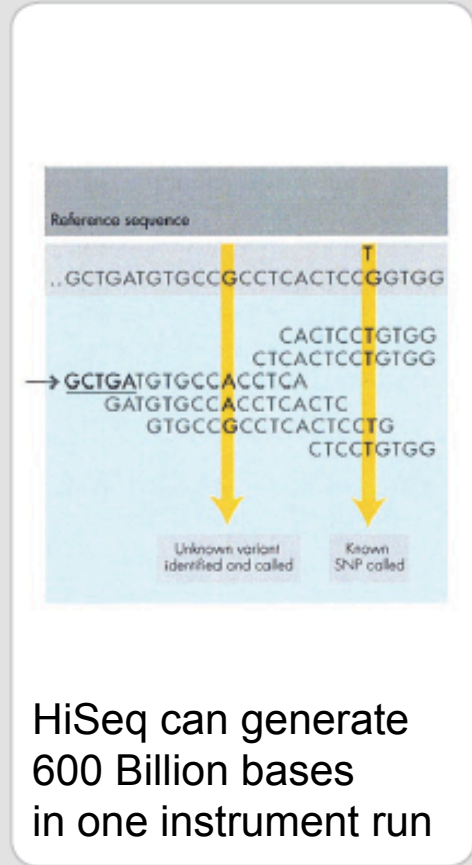
After laser excitation, the image is captured as before, and the identity of the second base is recorded.

## 11. SEQUENCING OVER MULTIPLE CHEMISTRY CYCLES



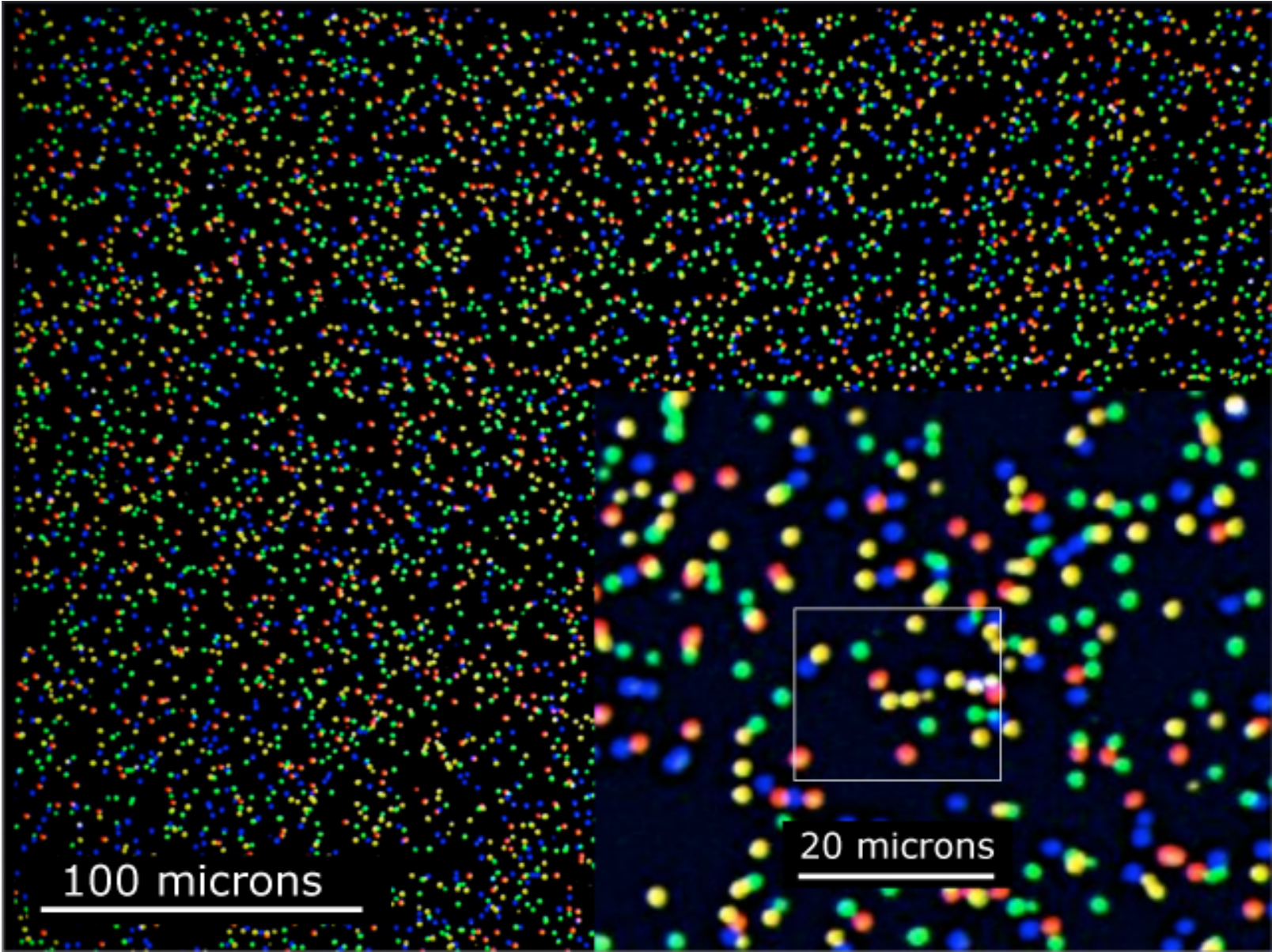
The sequencing cycles are repeated to determine the sequence of bases in a fragment, one base at a time.

## 12. ALIGN DATA



HiSeq can generate 600 Billion bases in one instrument run

The data are aligned and compared to a reference, and sequencing differences are identified.



# Sequencing Analysis Viewer

Run Folder: Y:\111208\_SN372\_0101\_AD0JRMACXX

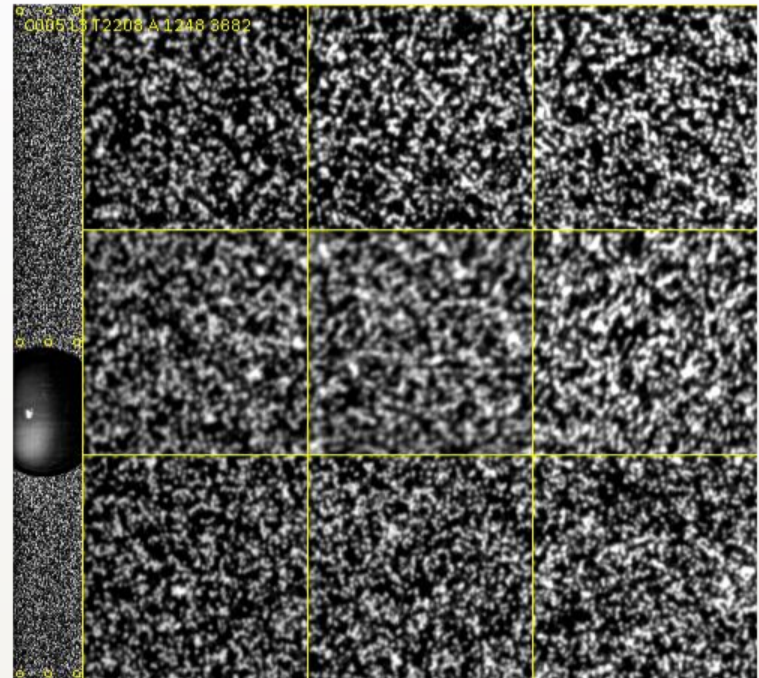
**Browse** **Refresh**

Analysis **Imaging** Summary Tile Status Controls

Cycle All Lane 3 Surface Bottom Swath Middle Section All

A  C  G  T

Index	Lane	Tile	Section	Cycle	Surface	Swath	Time	P90 A	P90 C	P90 G	P90 T
33147	3	2207	7	99	Bottom	Middle	12/13/201...	1077	2909	932	22
33148	3	2207	7	100	Bottom	Middle	12/13/201...	1871	2966	931	22
33149	3	2207	7	101	Bottom	Middle	12/13/201...	1840	2929	919	22
33150	3	2207	7	102	Bottom	Middle	12/13/201...	1822	2893	922	22
33151	3	2207	7	103	Bottom	Middle	12/14/201...	1805	2878	906	22
33152	3	2207	7	104	Bottom	Middle	12/14/201...	1802	2876	902	21
33153	3	2207	7	105	Bottom	Middle	12/14/201...	1785	2841	908	21
33154	3	2207	7	106	Bottom	Middle	12/14/201...	1756	2836	874	21
33155	3	2207	7	107	Bottom	Middle	12/14/201...	1749	2813	872	21
33156	3	2207	7	108	Bottom	Middle	12/14/201...	2498	3963	938	25
33913	3	2208	8	1	Bottom	Middle	12/08/201...	2928	4976	1861	35
33914	3	2208	8	2	Bottom	Middle	12/08/201...	3176	4792	1636	41
33915	3	2208	8	3	Bottom	Middle	12/08/201...	3163	4773	1679	36
33916	3	2208	8	4	Bottom	Middle	12/08/201...	3259	4788	1690	34
<b>33917</b>	<b>3</b>	<b>2208</b>	<b>8</b>	<b>5</b>	<b>Bottom</b>	<b>Middle</b>	<b>12/08/201...</b>	<b>2732</b>	<b>4112</b>	<b>1533</b>	<b>28</b>
33918	3	2208	8	6	Bottom	Middle	12/09/201...	3126	4605	1475	33
33919	3	2208	8	7	Bottom	Middle	12/09/201...	2712	4312	1404	36



Rows=41472 Disp=864 Sel=1 Filter

# Sequencing Analysis Viewer

Run Folder: Y:\111208\_SN372\_0101\_AD0JRMACXX

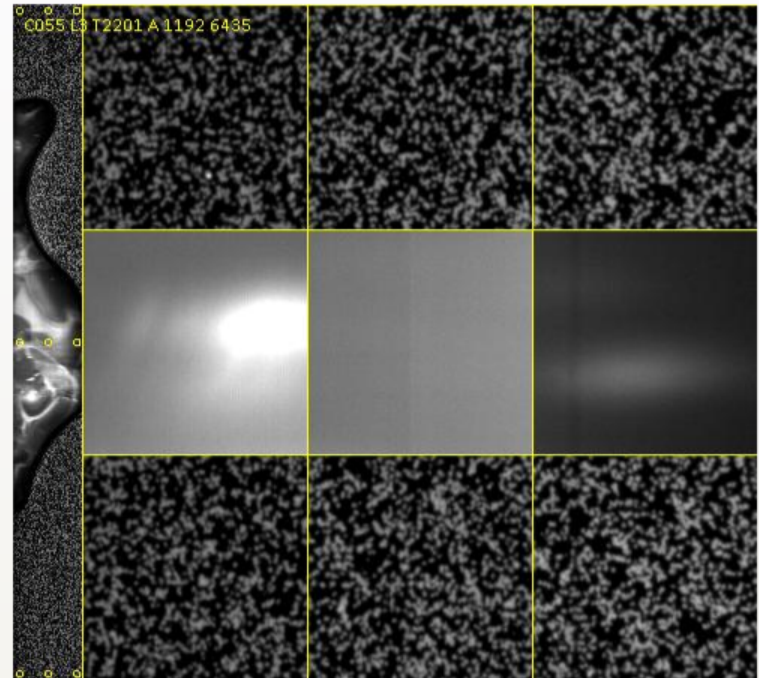
**Browse** **Refresh**

Analysis **Imaging** Summary Tile Status Controls

Cycle All Lane 3 Surface Bottom Swath Middle Section All

A  C  G  T

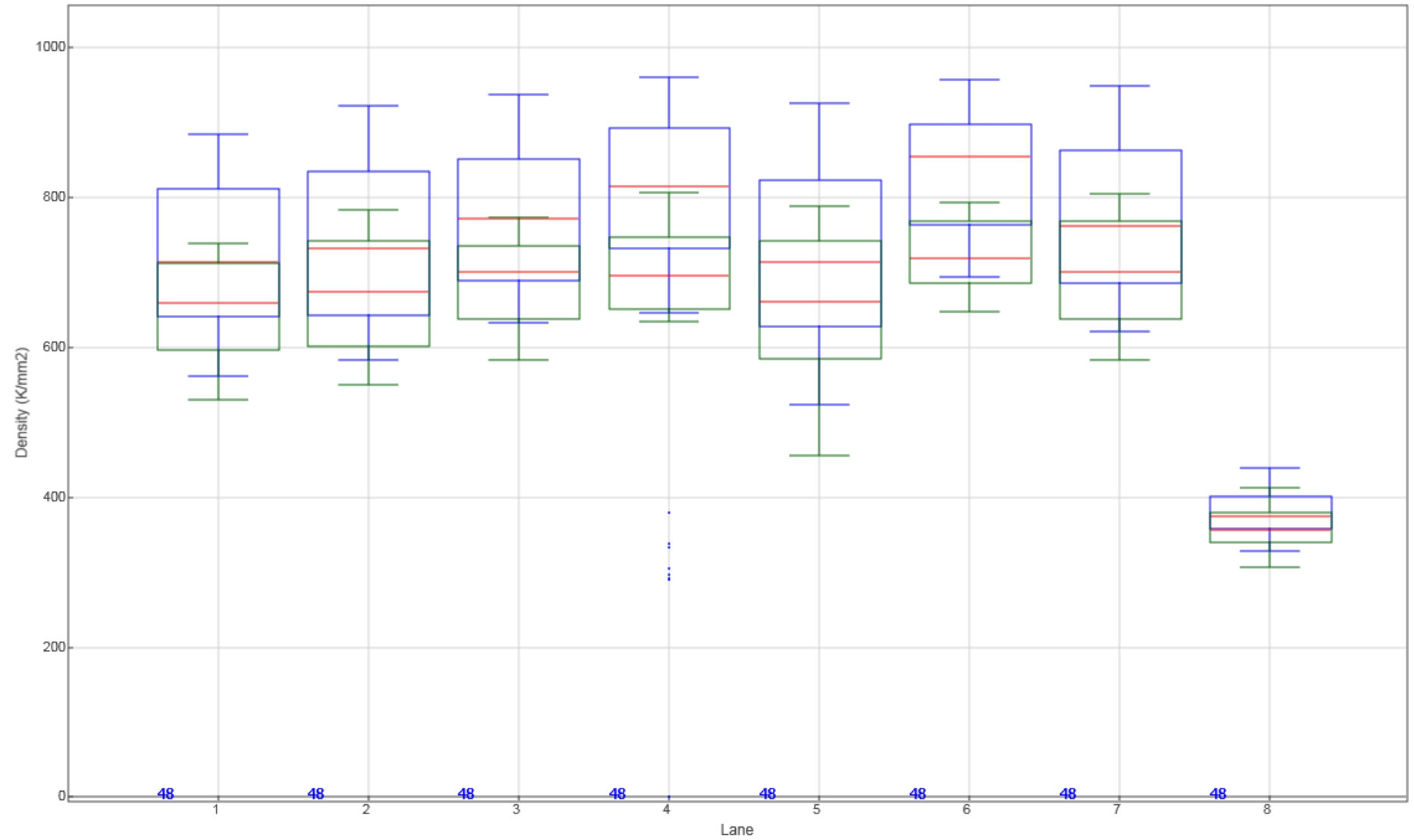
Index	Lane	Tile	Section	Cycle	Surface	Swath	Time	P90 A	P90 C	P90 G	P90 T
27905	3	2201	1	41	Bottom	Middle	12/10/201...	2368	3639	1103	26
27906	3	2201	1	42	Bottom	Middle	12/10/201...	2369	3824	1148	27
27907	3	2201	1	43	Bottom	Middle	12/10/201...	2361	3831	1136	27
27908	3	2201	1	44	Bottom	Middle	12/10/201...	2301	3759	1106	26
27909	3	2201	1	45	Bottom	Middle	12/10/201...	156	218	146	18
27910	3	2201	1	46	Bottom	Middle	12/10/201...	2298	3726	1124	26
27911	3	2201	1	47	Bottom	Middle	12/10/201...	2263	3673	1100	26
27912	3	2201	1	48	Bottom	Middle	12/10/201...	2211	3644	1065	25
27913	3	2201	1	49	Bottom	Middle	12/10/201...	2228	3657	1068	25
27914	3	2201	1	50	Bottom	Middle	12/10/201...	150	206	0	0
27915	3	2201	1	51	Bottom	Middle	12/10/201...	3761	5324	2061	43
27916	3	2201	1	52	Bottom	Middle	12/10/201...	3608	5397	1707	11
27917	3	2201	1	53	Bottom	Middle	12/10/201...	0	0	176	32
27918	3	2201	1	54	Bottom	Middle	12/10/201...	961	4845	167	44
27919	3	2201	1	55	Bottom	Middle	12/10/201...	2538	4054	1430	26
27920	3	2201	1	56	Bottom	Middle	12/10/201...	2813	4531	1552	34
27921	3	2201	1	57	Bottom	Middle	12/10/201...	300	566	0	0
27922	3	2201	1	58	Bottom	Middle	12/10/201...	0	0	0	0



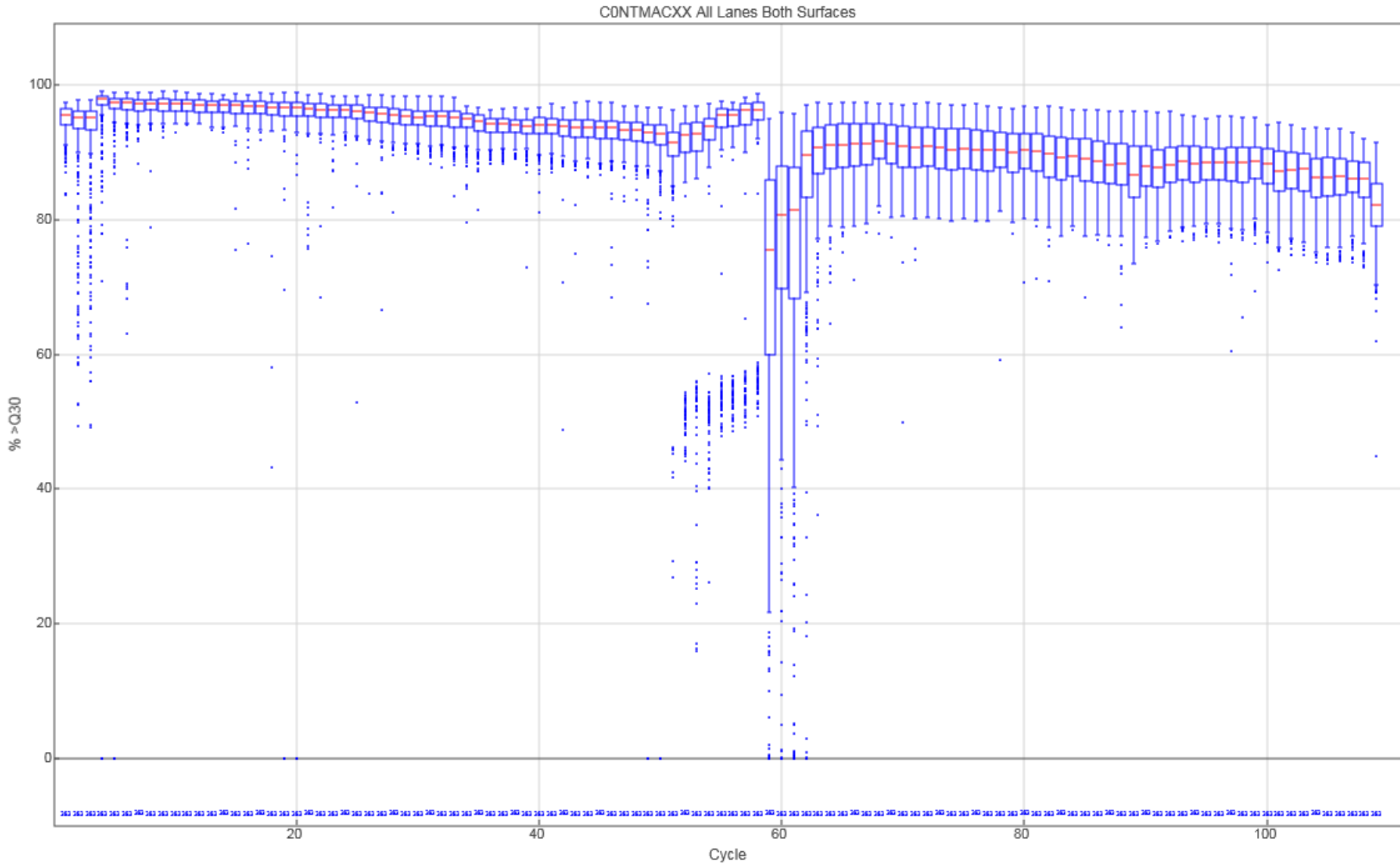
Rows=41472 Disp=864 Sel=1 Filter

# Cluster Density

CONTMACXX Both Surfaces

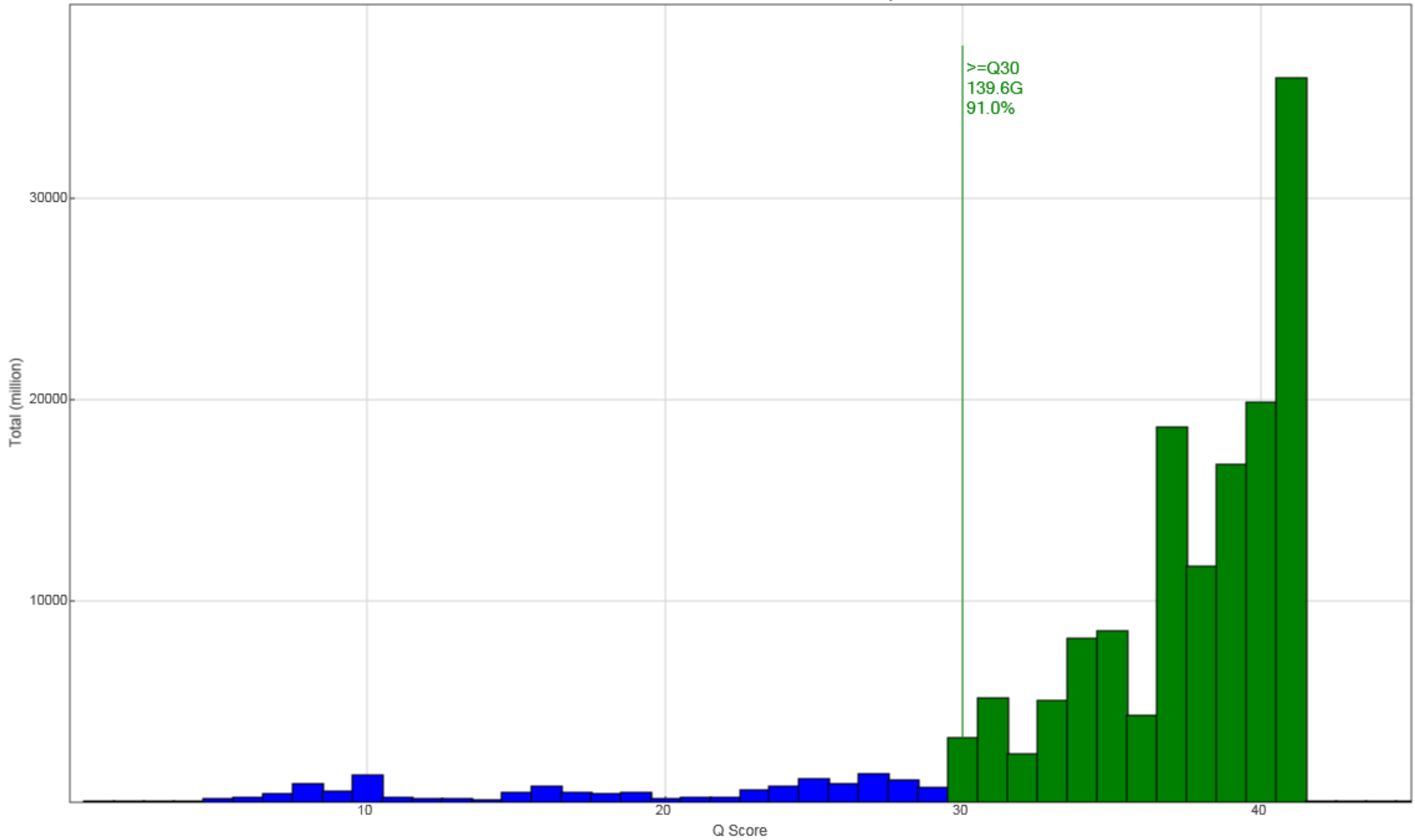


# Percent Q30 Scores per cycle for all lanes and both surfaces



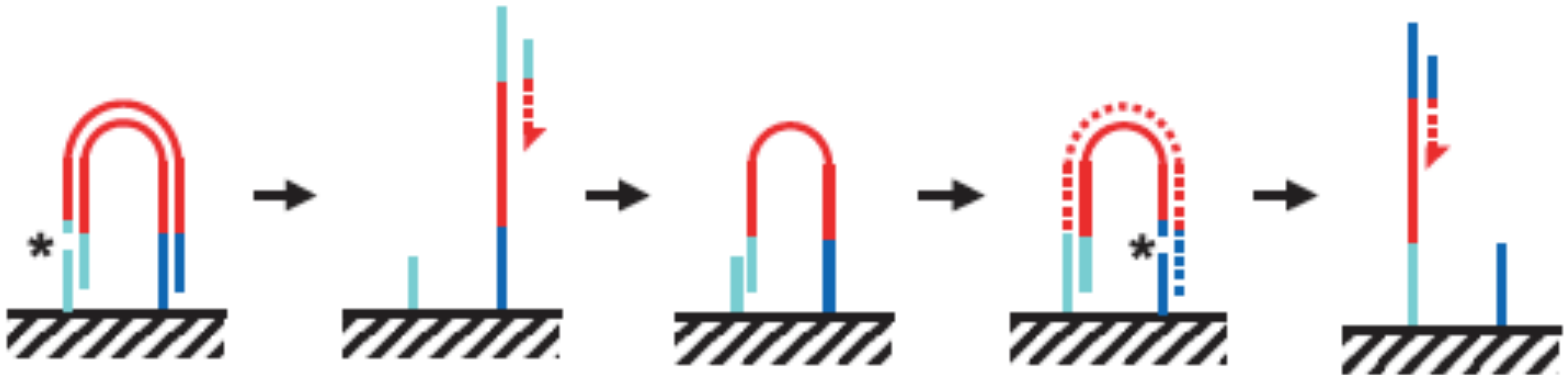
# Q Scores

CONTMACXX All Lanes Both Surf. All Reads All Cycles



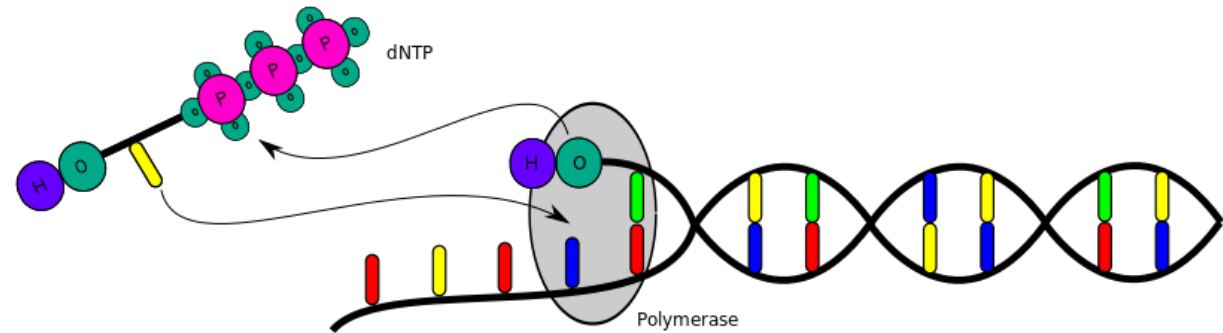
# Single End vs. Paired End Sequencing

c

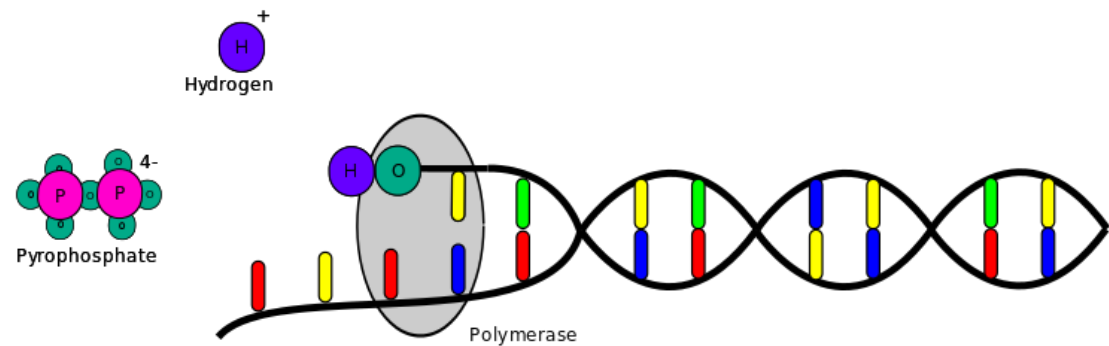




# Ion Semiconductor Sequencing (aka ion Torrent or Proton)

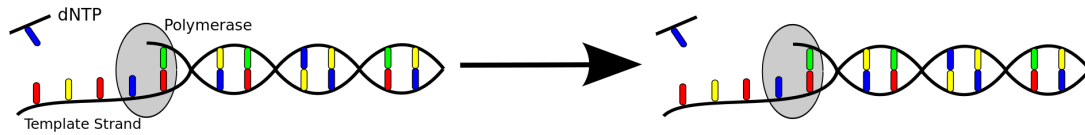


Polymerase integrates a nucleotide.

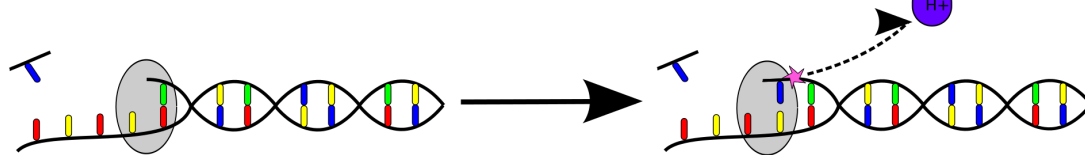


Hydrogen and pyrophosphate are released.

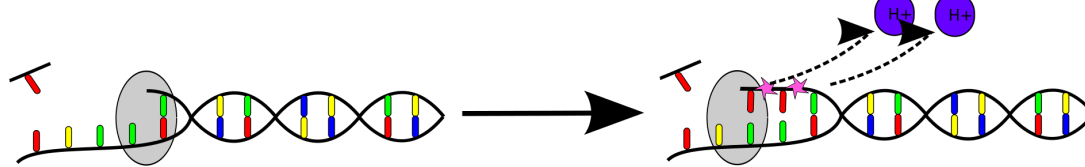
# Ion Semiconductor Sequencing



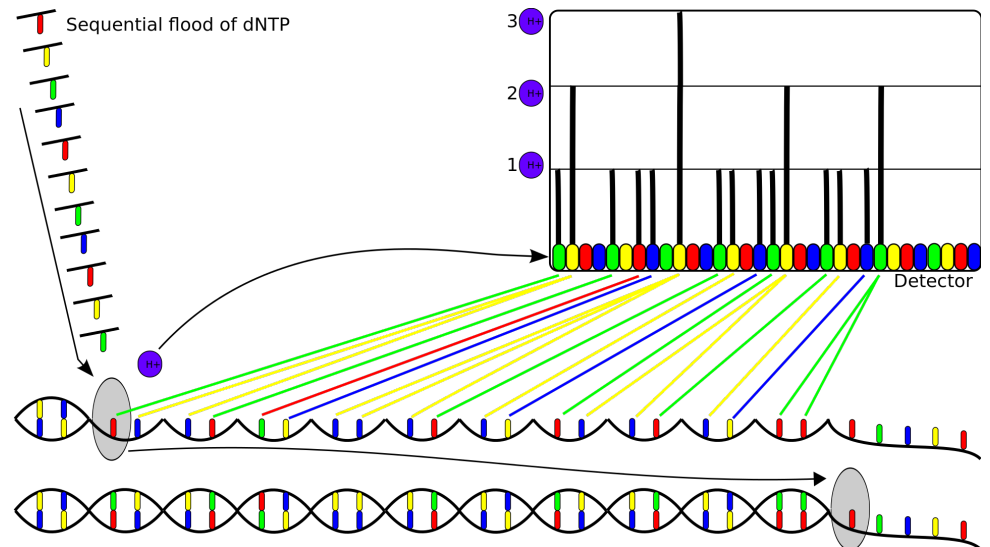
The nucleotide does not compliment the template - no release of hydrogen.



The nucleotide compliments the template - hydrogen is released.

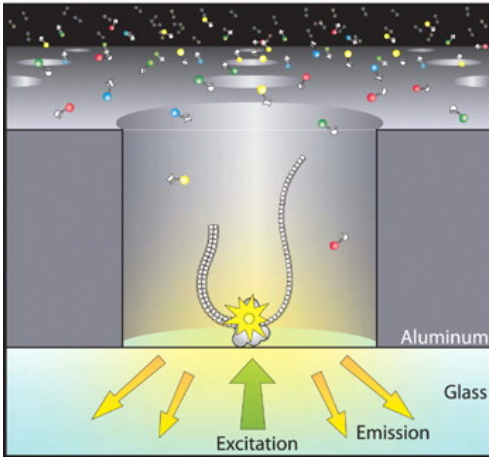


The nucleotide compliments several bases in a row - multiple hydrogen ions are released.

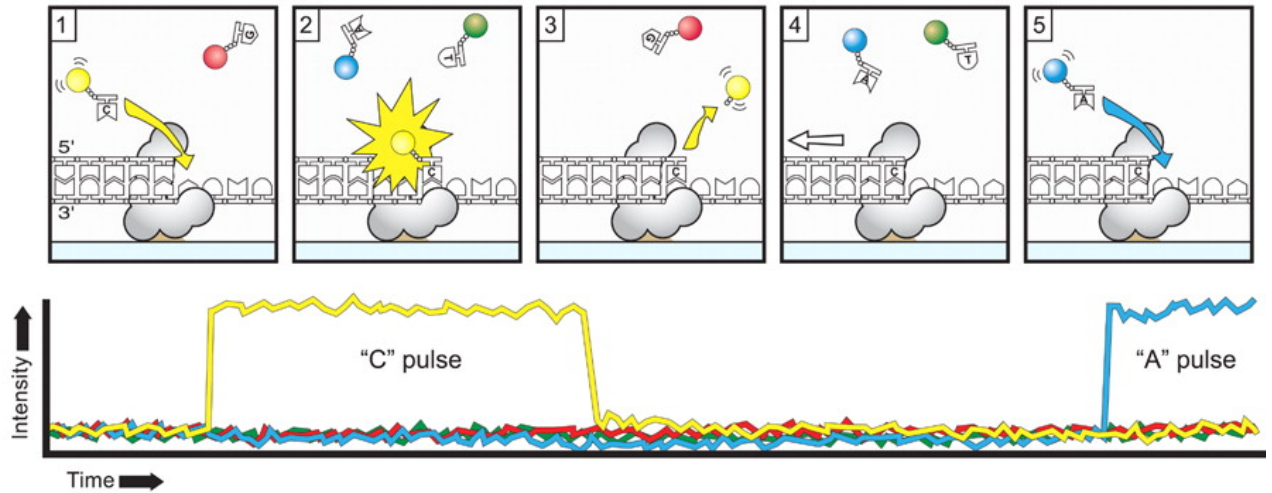


# Pacific Biosciences Technology

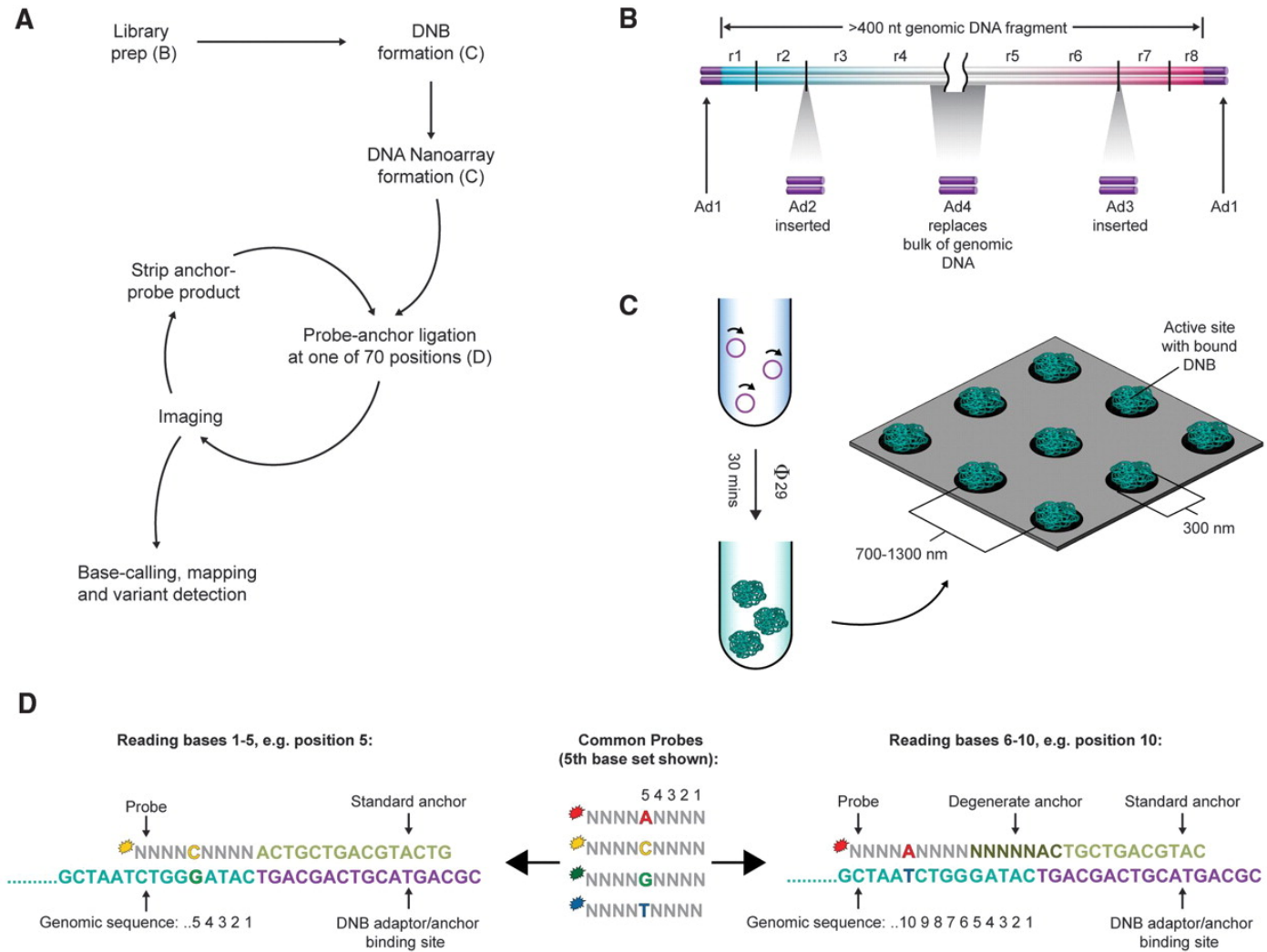
A



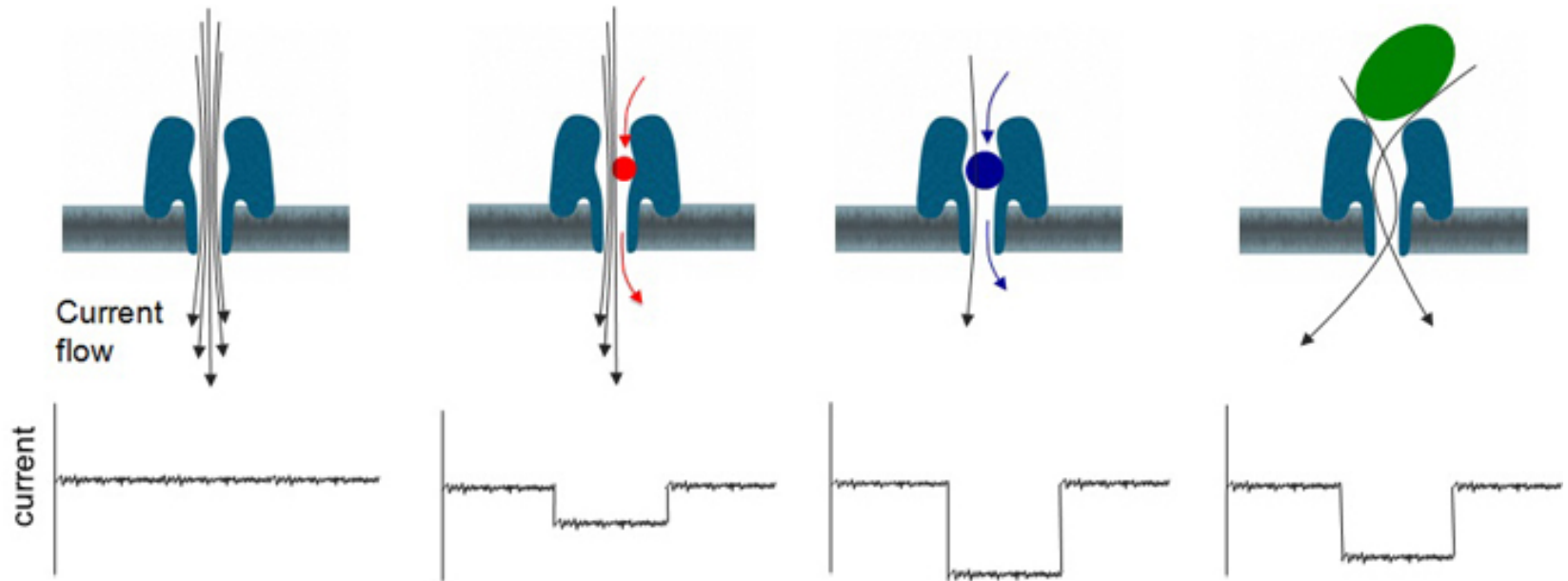
B



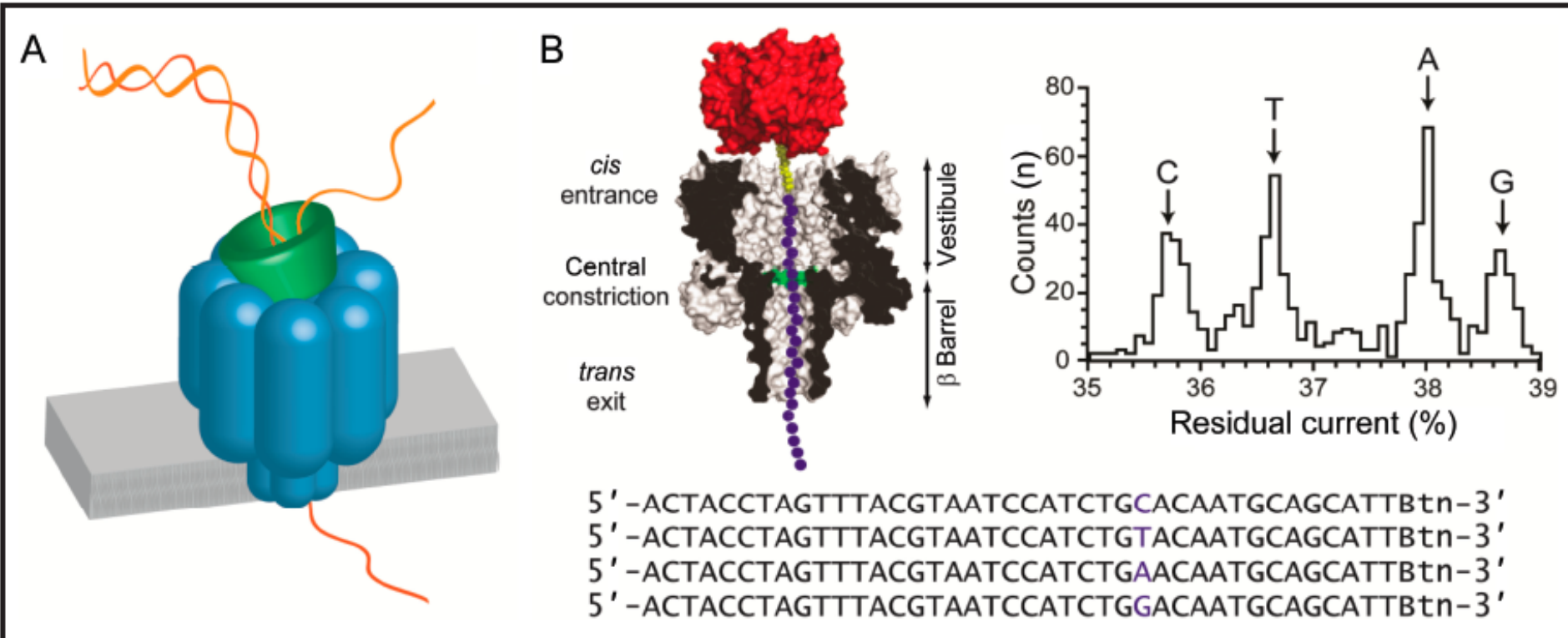
# Complete Genomics Technology



# Nanopores



# Oxford Nanopore



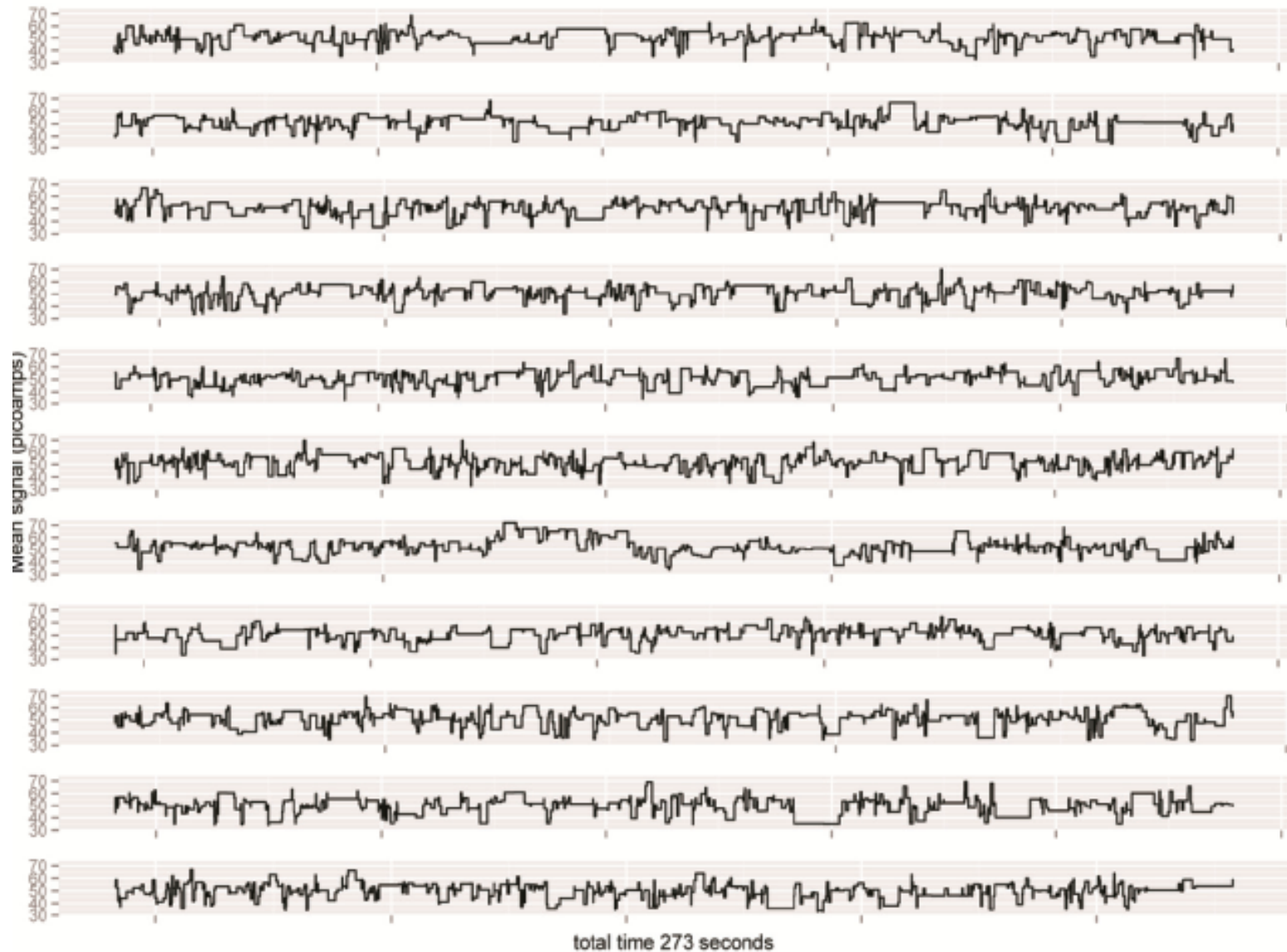
**Fig. 2. Nanopore strand sequencing.**

(A), Basis of nanopore sequencing. ssDNA is fed through an individual protein pore by an enzyme that handles dsDNA. The sequence is determined by analysis of fluctuations in the ionic current. (B), Early base identification experiments. ssDNAs were suspended in an  $\alpha$ HL pore by attachment to streptavidin to mimic the ratcheting motion of the enzyme. The bases G, A, T, and C in a DNA hetero-oligomer each gave a different residual ionic current. Adapted with permission from Stoddart et al. (25).

# The MinION

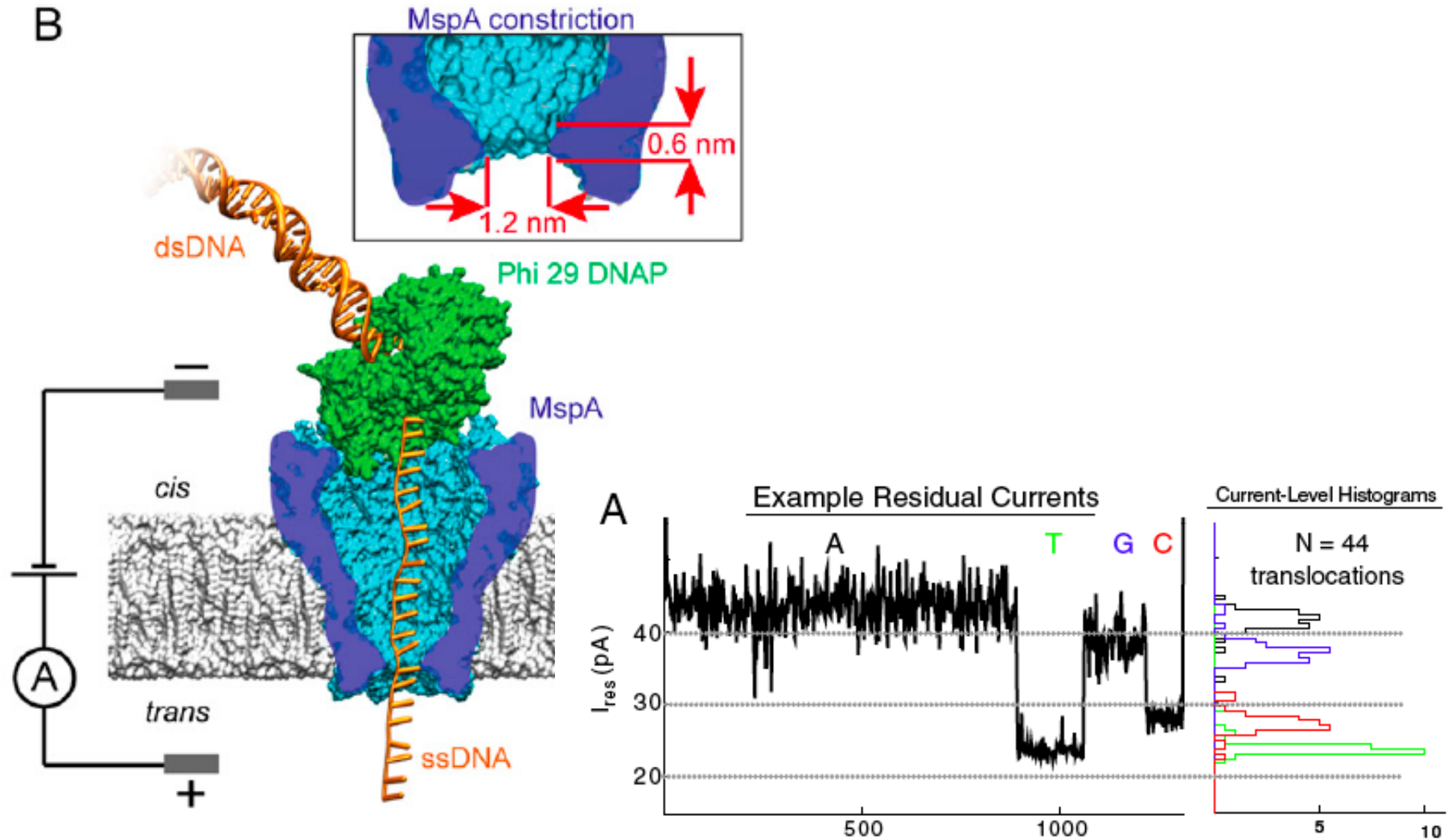


# Sequence from the minION





# MspA Nanopore



# Sequencing DNA

# Human Whole Genome Sequencing

- Initial Ref Sequence \$300 million and took about a decade. (Draft reported in 2001)
- Humans sequenced
  - Craig Venter
  - James Watson
  - Yoruban from HapMap
  - Korean (35 individuals published in 2014)
  - Han Chinese
- Broad has released 60K Human Exomes
  - [exac.broadinstitute.org](http://exac.broadinstitute.org)
- Broad has “tweeted” it has completed 10K whole human genomes with the Illumina X Ten
- HiSeq2500 High Output--human genome can be sequenced for about \$5,000 at an average read depth of 30X in 10 days
- HiSeq2500 Rapid Run— human genome can be sequenced in about 2 days to 30X coverage for ~\$4,000.

# DNA Sequencing with Next-Generation Technologies

## A Draft Sequence of the Neandertal Genome

Richard E. Green,<sup>1\*</sup>†‡ Johannes Krause,<sup>1†§</sup> Adrian W. Briggs,<sup>1†§</sup> Tomislav Maricic,<sup>1†§</sup> Udo Stenzel,<sup>1†§</sup> Martin Kircher,<sup>1†§</sup> Nick Patterson,<sup>2†§</sup> Heng Li,<sup>2†</sup> Weiwei Zhai,<sup>3†||</sup> Markus Hsi-Yang Fritz,<sup>4†</sup> Nancy F. Hansen,<sup>5†</sup> Eric Y. Durand,<sup>3†</sup> Anna-Sapfo Malaspinas,<sup>3†</sup> Jeffrey M. Rosenfeld,<sup>6†</sup> Ronan O. Quinlan,<sup>7,13†</sup> Xiaohu Wu,<sup>7†</sup> David A. Preiss,<sup>1†</sup> et al.



Hernán  
Barbar  
Eric S.  
Christi  
Vladim  
Javier F  
Daniel  
Janet K

## Sequencing the nuclear genome of the extinct woolly mammoth

Webb Miller<sup>1</sup>, Daniela I. I.  
Michael D. Packard<sup>1</sup>, Fan  
Kerstin Lindblad-Toh<sup>5</sup>, Eri  
Sharon Sheridan<sup>7</sup>, Tom P

## Genetic history of an archaic hominin group from Denisova Cave in Siberia

David Reich<sup>1,2\*</sup>, Richard E. Green<sup>3,4\*</sup>, Martin Kircher<sup>3\*</sup>, Johannes Krause<sup>3,5\*</sup>, Nick Patterson<sup>2\*</sup>, Eric Y. Durand<sup>6\*</sup>, Rence Viola<sup>3,7\*</sup>

Adrian W. Briggs<sup>1,3</sup>, U  
Can Alkan<sup>10</sup>, Qiaomei  
Michael Richards<sup>7,13</sup>,  
Montgomery Slatkin<sup>6</sup>

## Insights into hominid evolution from the gorilla genome sequence



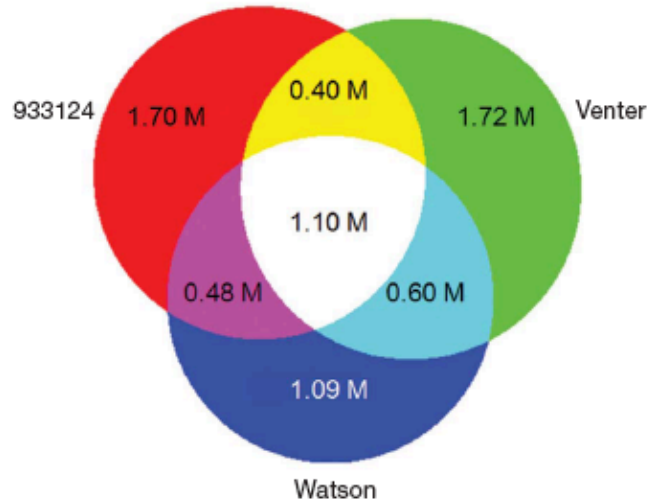
Aylwyn Scally<sup>1</sup>, Julien Y. Duthéil<sup>2†</sup>, LaDeana W. Hillier<sup>3</sup>, Gregory E. Jordan<sup>4</sup>, Ian Goodhead<sup>1†</sup>, Javier Herrero<sup>4</sup>, Asger Hobolth<sup>2</sup>, Tuuli Lappalainen<sup>5</sup>, Thomas Mailund<sup>2</sup>, Tomas Marques-Bonet<sup>3,6,7</sup>, Shane McCarthy<sup>1</sup>, Stephen H. Montgomery<sup>8</sup>, Petra C. Schwalie<sup>4</sup>, Y. Amy Tang<sup>1</sup>, Michelle C. Ward<sup>9,10</sup>, Yali Xue<sup>1</sup>, Bryndis Yngvadottir<sup>1†</sup>, Can Alkan<sup>3,11</sup>, Lars N. Andersen<sup>2</sup>, Qasim Ayub<sup>1</sup>, Edward V. Ball<sup>12</sup>, Kathryn Beal<sup>4</sup>, Brenda J. Bradley<sup>8,13</sup>, Yuan Chen<sup>1</sup>, Chris M. Clee<sup>1</sup>, Stephen Fitzgerald<sup>4</sup>, Tina A. Graves<sup>14</sup>, Yong Gu<sup>1</sup>, Paul Heath<sup>1</sup>, Andreas Heger<sup>15</sup>, Emre Karakoc<sup>3</sup>, Anja Kolb-Kokocinski<sup>1</sup>, Gavin K. Laird<sup>1</sup>, Gerton Lunter<sup>16</sup>, Stephen Meader<sup>15</sup>, Matthew Mort<sup>12</sup>, James C. Mullikin<sup>17</sup>, Kasper Munch<sup>2</sup>, Timothy D. O'Connor<sup>8</sup>, Andrew D. Phillips<sup>12</sup>, Javier Prado-Martinez<sup>9</sup>, Anthony S. Rogers<sup>1†</sup>, Saba Sajadian<sup>3</sup>, Dominic Schmidt<sup>9,10</sup>, Katy Shaw<sup>12</sup>, Jared T. Simpson<sup>1</sup>, Peter D. Stenson<sup>12</sup>, Daniel J. Turner<sup>1†</sup>, Linda Vigilant<sup>18</sup>, Albert J. Vilella<sup>4</sup>, Weldon Whitener<sup>1</sup>, Baoli Zhu<sup>19†</sup>, David N. Cooper<sup>12</sup>, Pieter de Jong<sup>19</sup>, Emmanouil T. Dermitzakis<sup>5</sup>, Evan E. Eichler<sup>3,11</sup>, Paul Flicek<sup>4</sup>, Nick Goldman<sup>4</sup>, Nicholas I. Mundy<sup>8</sup>, Zemin Ning<sup>1</sup>, Duncan T. Odom<sup>1,9,10</sup>, Chris P. Ponting<sup>15</sup>, Michael A. Quail<sup>1</sup>, Oliver A. Ryder<sup>20</sup>, Stephen M. Searle<sup>1</sup>, Wesley C. Warren<sup>14</sup>, Richard K. Wilson<sup>14</sup>, Mikkel H. Schierup<sup>2</sup>, Jane Rogers<sup>1†</sup>, Chris Tyler-Smith<sup>1</sup> & Richard Durbin<sup>1</sup>

# Applications

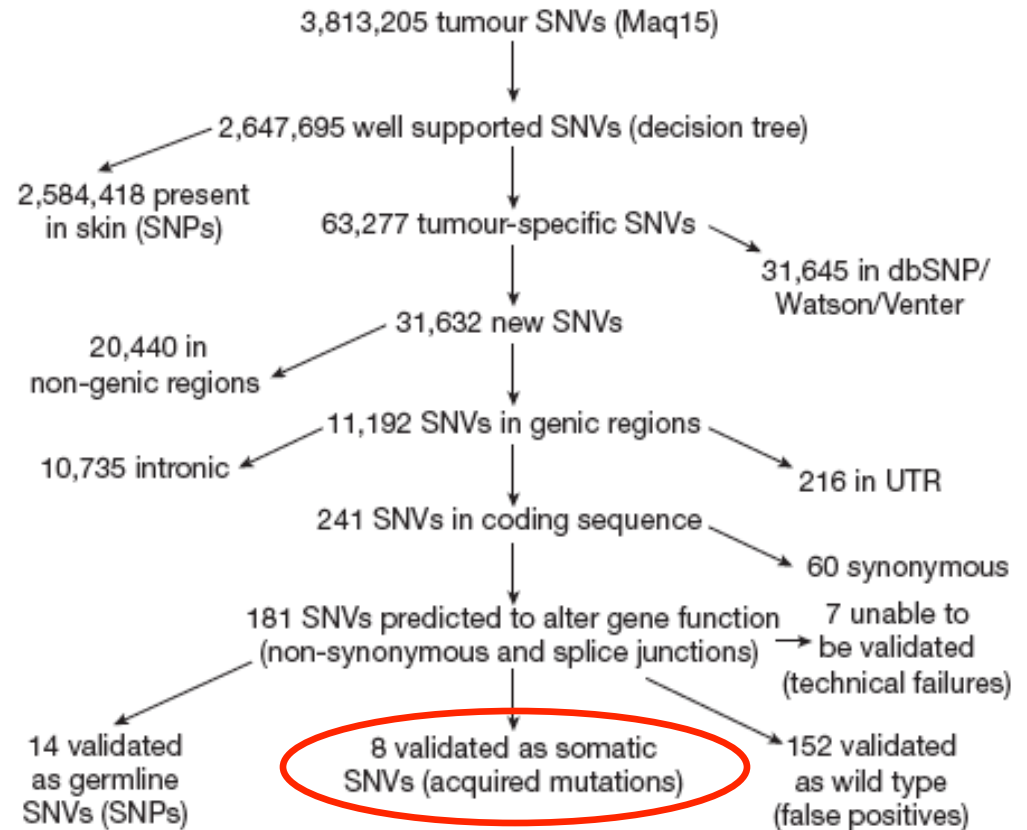
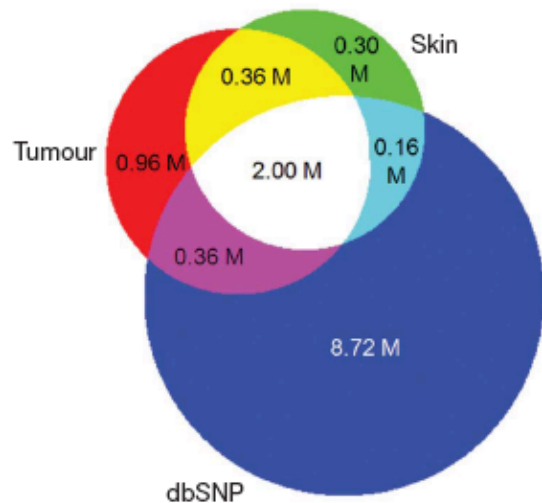
- Whole Genome Sequencing
- Exome Sequencing
- Targeted Genomic Sequencing
- Chromatin-IP-Sequencing
- DNase I Hypersensitivity Sequencing
- Methyl-Seq (RRBS, MeDIP, etc)
- Microbiome Sequencing
- Metagenomics

# AML: Comparisons

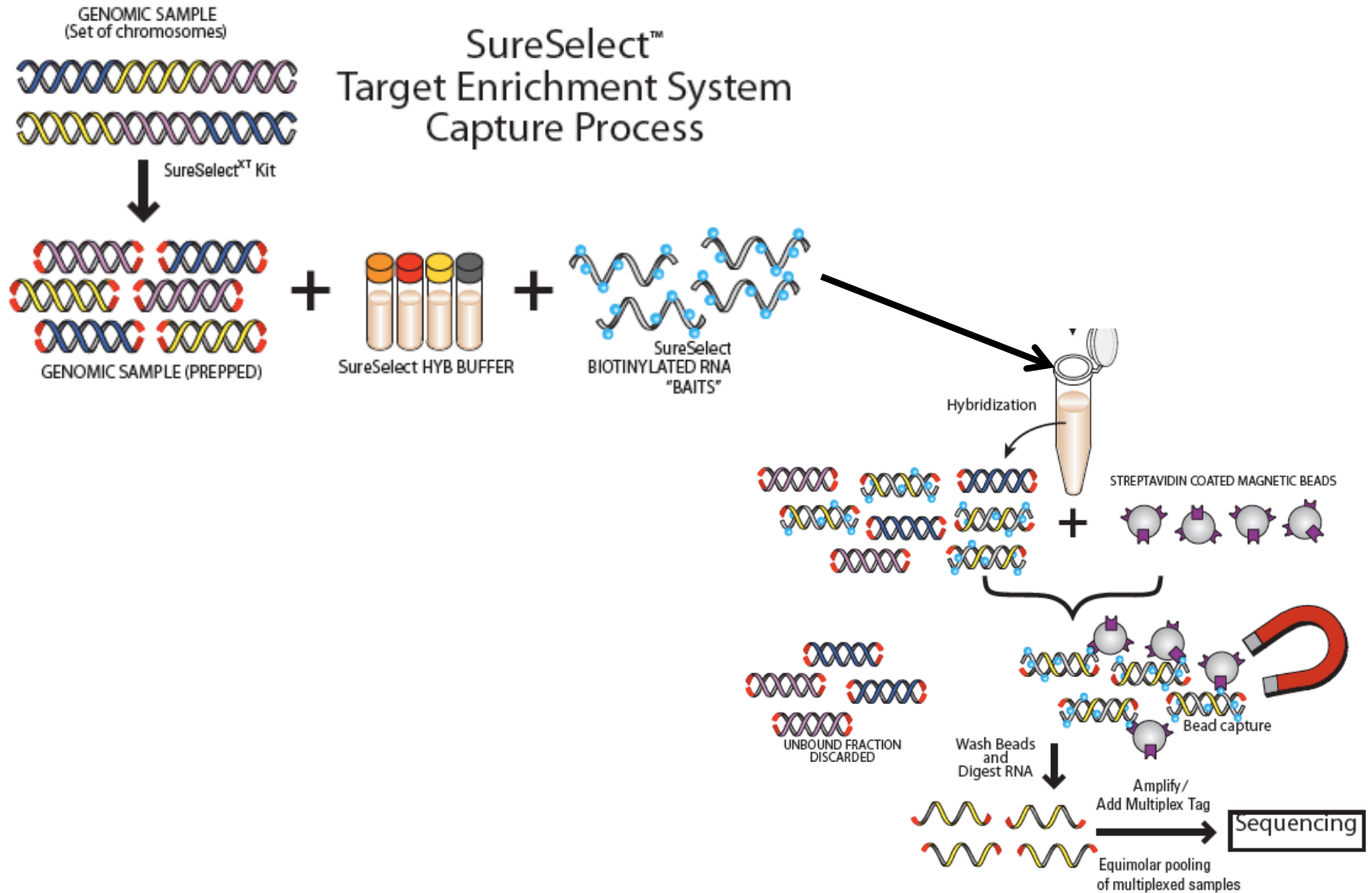
**a**



**b**



# SureSelect Exome Capture



# Disease Genes Discovered by Direct Whole Exome Sequencing\*

Gene Identified	Disease/Syndrome	Reference
MYH3	Freeman-Sheldon Syndrome	Ng SB, et al. 2009. Nature 462
SLC26A3	Bartter Syndrome	Choi M, et al. 2009 PNAS 106(45)
DHODH	Miller Syndrome	Ng SB, et al. 2010 Nat Genet 42(1).
FLVCR2	Fowler Syndrome	Lalonde, E. et al. 2010 Hum Mutat 31(8).
FLNA	Terminal Osseous Dysplasia (TOD)	Sun Y., et al. 2010 Am J. Hum Genet 87(1).
GPSM2	Nonsyndromic Hearing Loss (DFNB82)	Walsh, T. et al. 2010 Am J. Hum Genet 87(1).
HSD17B4	Perrault Syndrome/DBP	Pierce SB, et al. 2010 Am J. Hum Genet 87(2).
MLL2	Kabuki Syndrome	Ng SB, et al. 2010 Nat Genet 42(9).
ABCG5	Hypercholesterolemia	Rios J., et al. 2010 Hum Mol Genet 19(22).
WDR62	Brain Malformations	Bilguvar K, et al. 2010 Nature 467(7312).
PIGV	Hyperphosphatasia Mental Retardation (HPMR)	Krawitz PM, et al. 2010 Nat Genet 42(10)
WDR35	Sensenbrenner Syndrome	Gilissen C, et al. 2010Am J Hum Genet 87(3).
SDCCAG8	Nephromophthisis-related Ciliopathies	Otto EA, et al. 2010 Nat Genet 42(10).
STIM1	Kaposi Sarcoma	Byn M, et al. 2010 J Exp Med 207(11).
SCARF2	Van Den Ende-Gupta Syndrome	Anastasio N. et al. 2010 Am J Hum Genet 87(4).
C20orf54	Brown-Vialetto-Van Laere Syndrome	Green P, et al. 2010 Am J Hum Genet 86(3).
MASP1	Carnevale, Malpuech, OSA and Michels Syndromes	Sirmaci A, at al. 2010 Am J Hum Genet 87(5).
ABCC8	Neonatal Diabetes Mellitus	Bonnefond A, et al. 2010 PLoS One 5(10).
BAP-1	Metastasizing Uveal Melanomas	Harbour JW, et al. 2010 Science Nov 4 Epub.
ACAD9	Complex I Deficiency	Haack TB, et al. 2010 Nat Genet Nov 7 Epub.
DYNC1H1	Mental Retardation	Vissers LELM, et al. 2010 Nat Genet 10.1038/ng.712
RAB39A	Mental Retardation	Vissers LELM, et al. 2010 Nat Genet 10.1038/ng.712
YY1	Mental Retardation	Vissers LELM, et al. 2010 Nat Genet 10.1038/ng.712
DEAF1	Mental Retardation	Vissers LELM, et al. 2010 Nat Genet 10.1038/ng.712

\*As of 23 Nov. 2010



# Targeted Re-sequencing

The ability to capture specific sequences in the genome

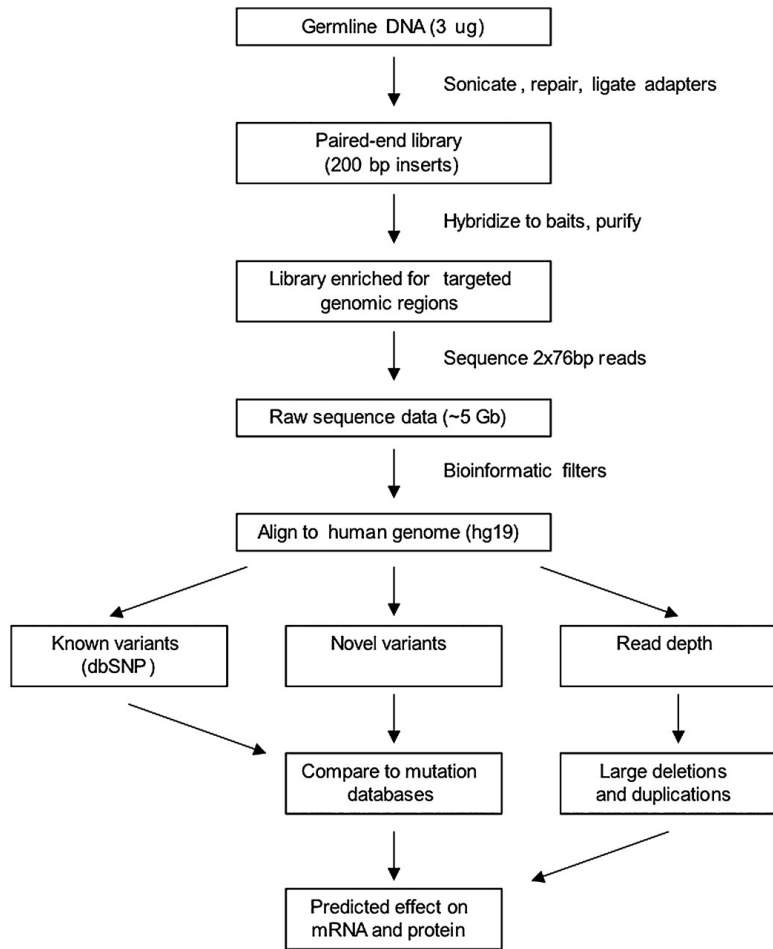
Long range PCR

Multiplex PCR strategies

Solution capture on Biotin labeled oligos

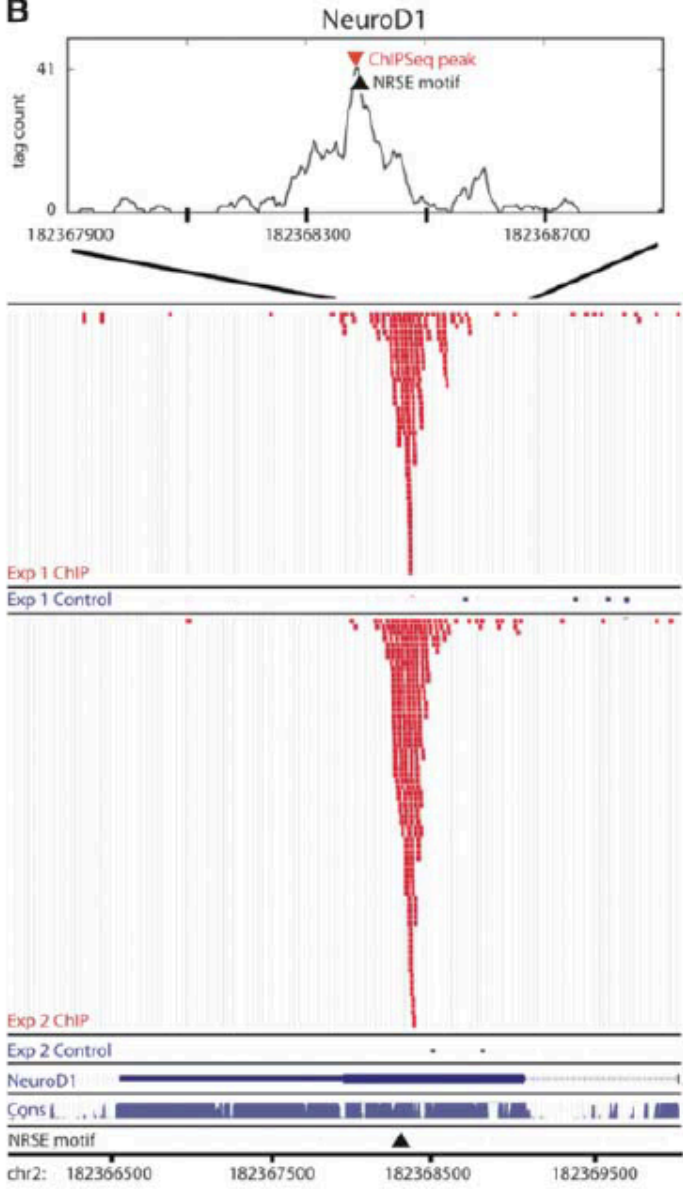
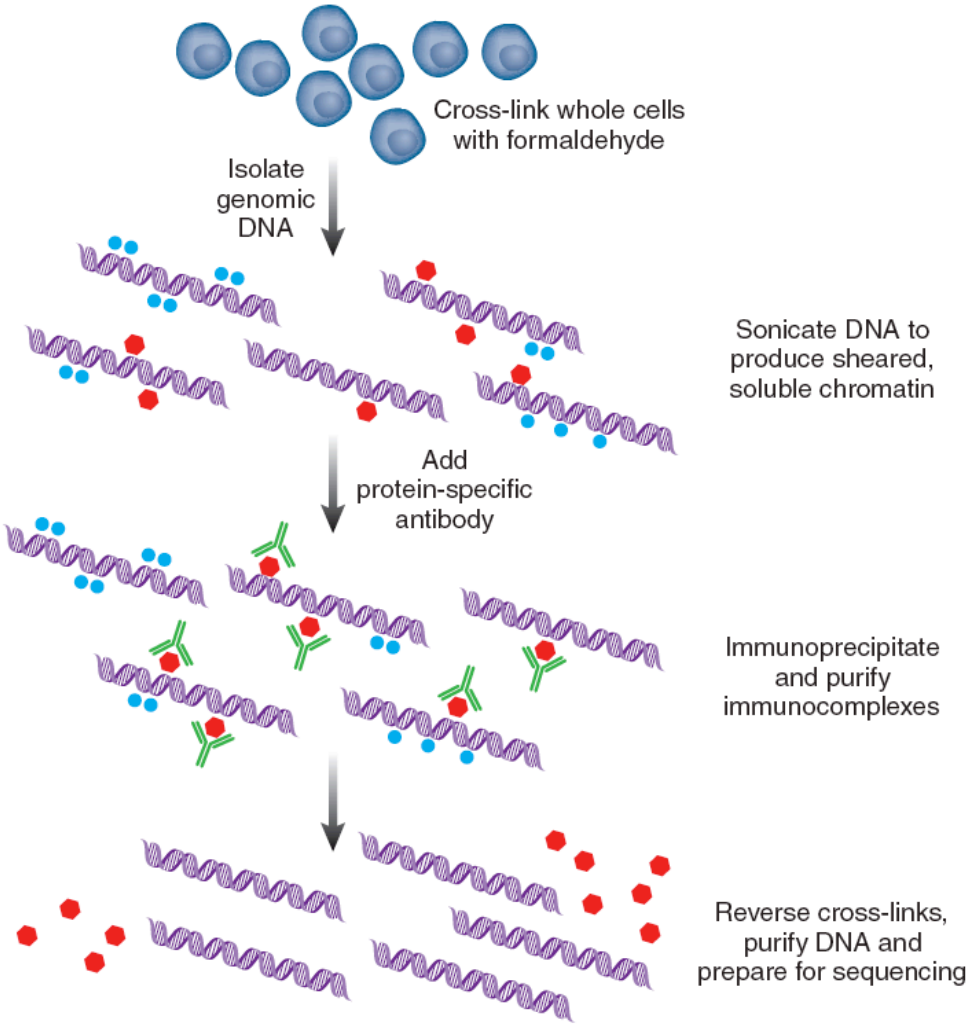
HaloPlex

# Genomic Capture of Breast Cancer Relevant Genes Followed by Next-Gen Sequencing.



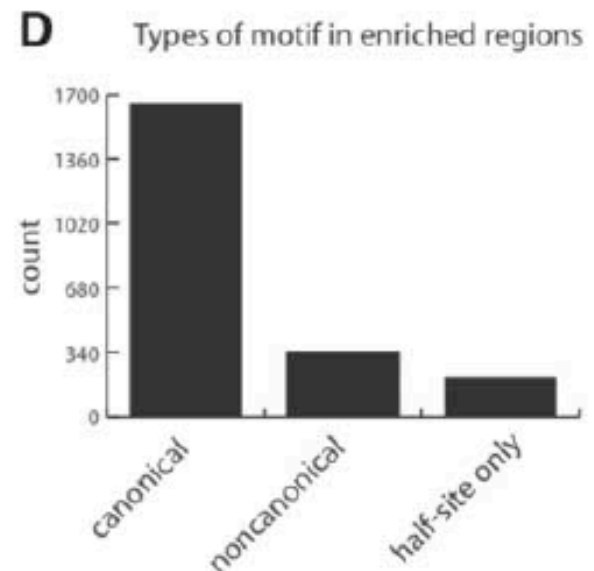
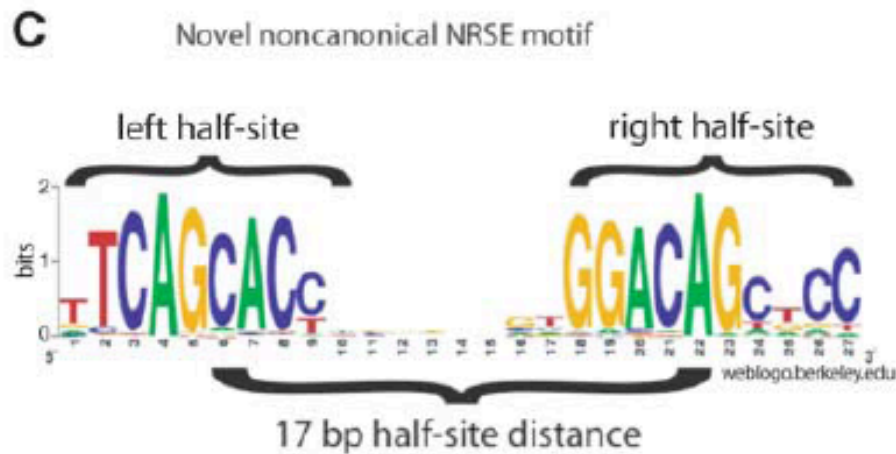
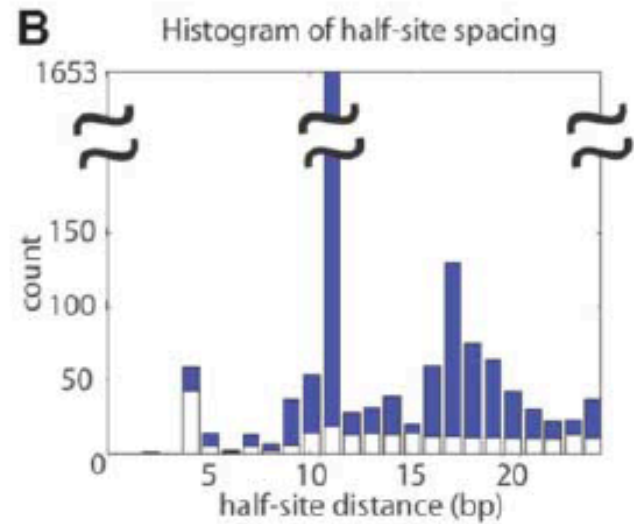
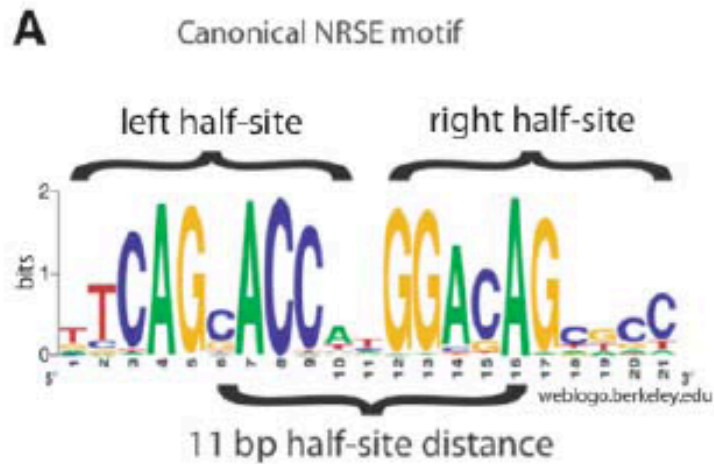
Gene	Chromosome	Start	End
BRCA1	17	41,186,313	41,347,712
BRCA2	13	32,879,617	32,983,809
CHEK2	22	29,073,731	29,147,822
PALB2	16	23,604,483	23,662,678
BRIP1	17	59,759,985	59,940,755
p53	17	7,561,720	7,600,863
PTEN	10	89,613,195	89,738,532
STK11	19	1,195,798	1,238,434
CDH1	16	68,761,195	68,879,444
ATM	11	108,083,559	108,249,826
BARD1	2	215,583,275	215,684,428
MLH1	3	37,024,979	37,102,337
MRE11	11	94,140,467	94,237,040
MSH2	2	47,620,263	47,720,360
MSH6	2	48,000,221	48,044,092
MUTYH	1	45,784,914	45,816,142
NBN	8	90,935,565	91,006,899
PMS1	2	190,638,811	190,752,355
PMS2	7	6,002,870	6,058,737
RAD50	5	131,882,630	131,989,595
RAD51C	17	56,759,963	56,821,692

# ChIP-Seq



**Figure 1** | Workflow of Chip-seq. DNA and proteins are cross-linked and purified; then bound DNA is analyzed by massively parallel short-read sequencing.

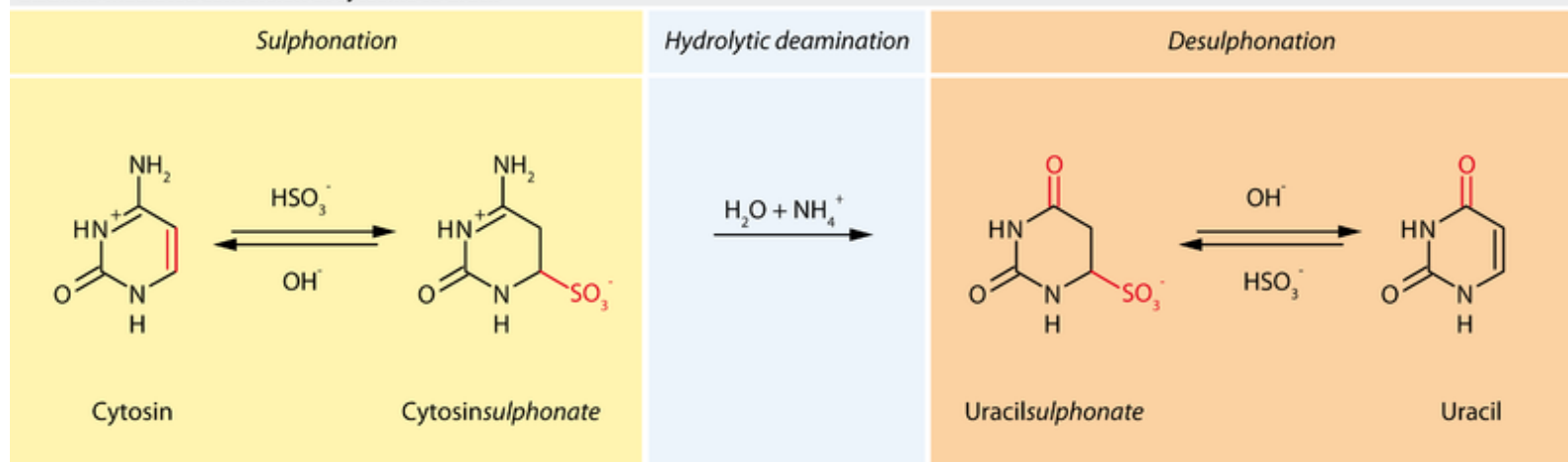
# ChIP-Seq



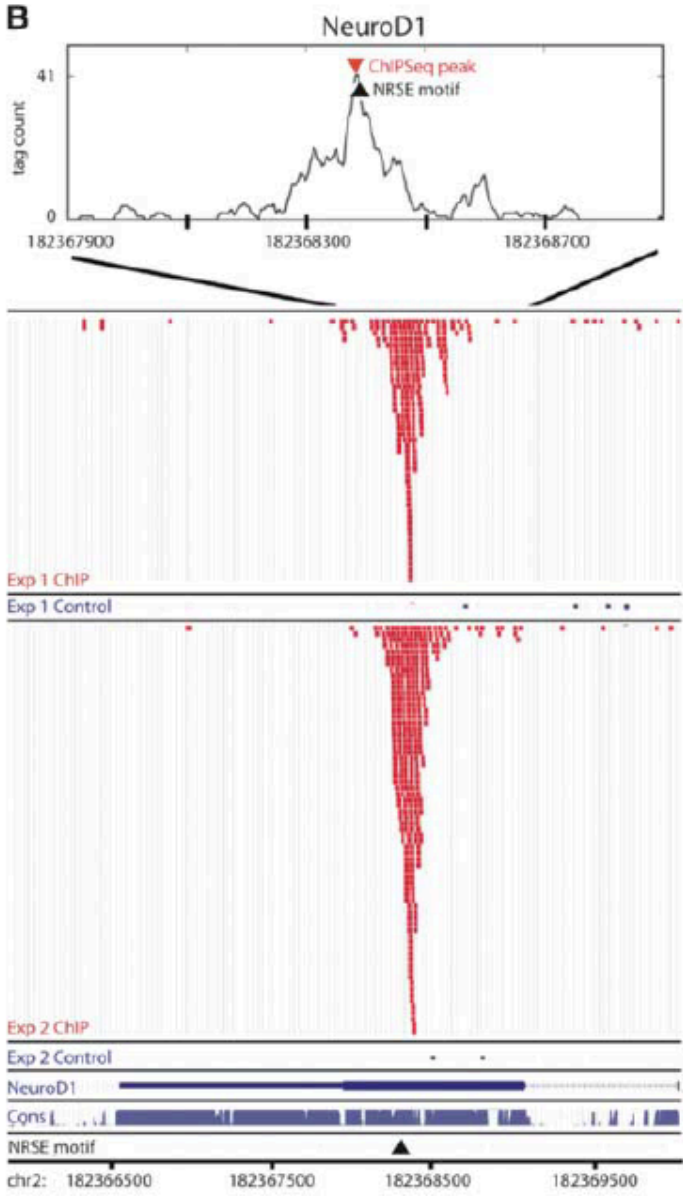
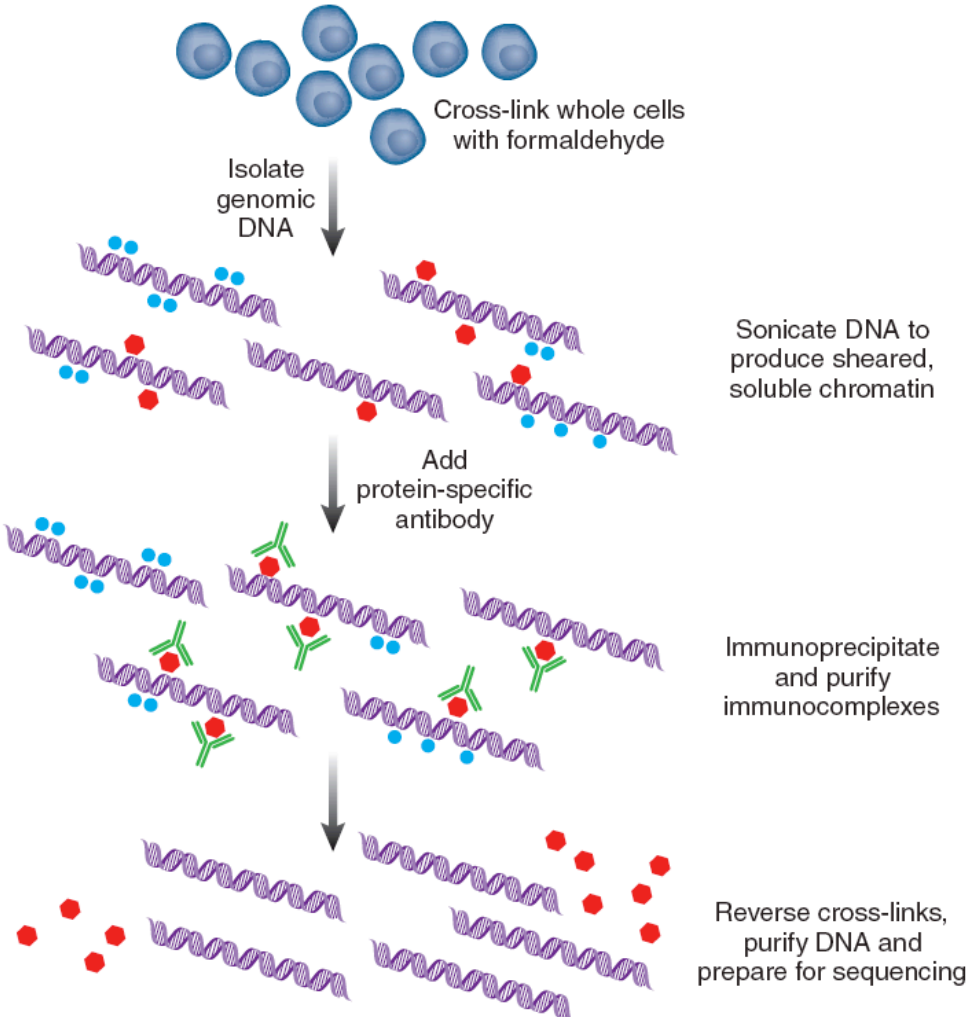
# Methylation profiling

- Whole genome bisulfite sequencing
- MeDIP (Methylated DNA-IP)
- Reduced Representational Bisulfite Sequencing
- Specific Capture methods

Bisulfite-mediated conversion of cytosine to uracil

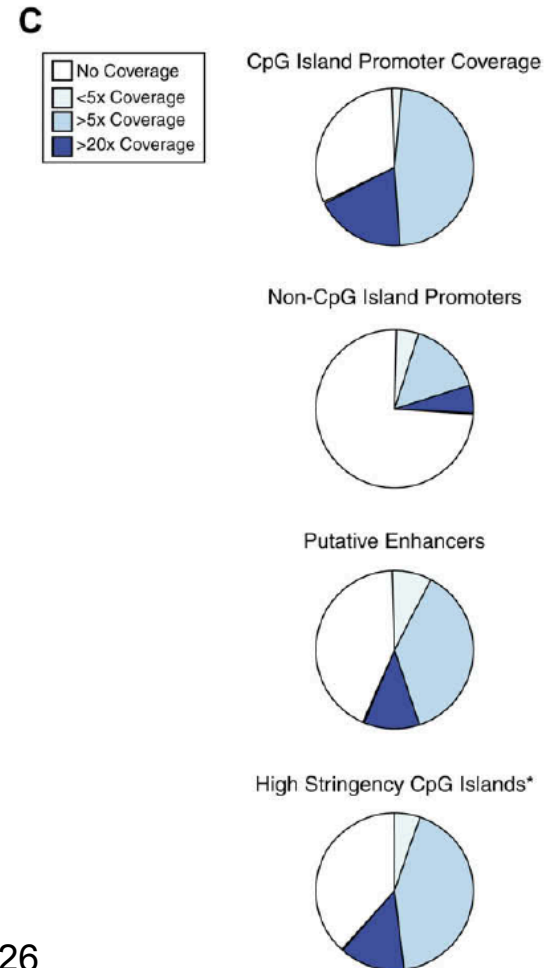
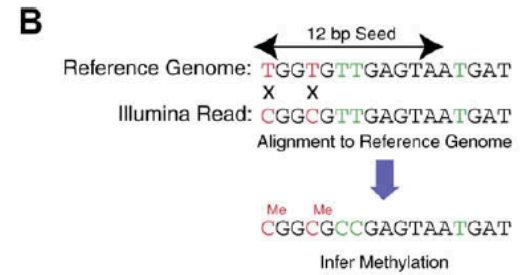
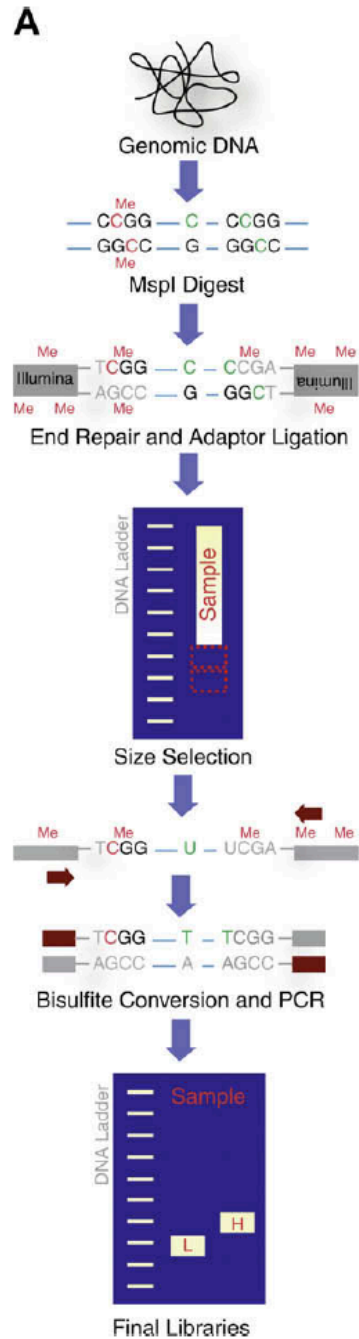


# MeDIP-Seq



**Figure 1** | Workflow of Chip-seq. DNA and proteins are cross-linked and purified; then bound DNA is analyzed by massively parallel short-read sequencing.

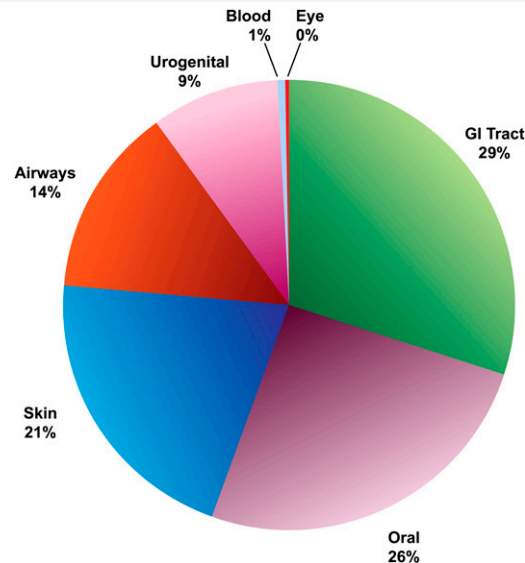
# RRBS



# The NIH Human Microbiome Project

## The NIH HMP Working Group<sup>1</sup>

The Human Microbiome Project (HMP), funded as an initiative of the NIH Roadmap for Biomedical Research (<http://nihroadmap.nih.gov>), is a multi-component community resource. The goals of the HMP are: (1) to take advantage of new, high-throughput technologies to characterize the human microbiome more fully by studying samples from multiple body sites from each of at least 250 “normal” volunteers; (2) to determine whether there are associations between changes in the microbiome and health/disease by studying several different medical conditions; and (3) to provide both a standardized data resource and new technological approaches to enable such studies to be undertaken broadly in the scientific community. The ethical, legal, and social implications of such research are being systematically studied as well. **The ultimate objective of the HMP is to demonstrate that there are opportunities to improve human health through monitoring or manipulation of the human microbiome. The history and implementation of this new program are described here.**



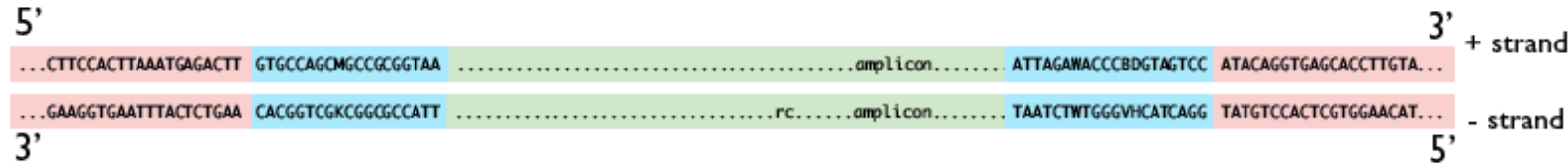
**Figure 3.** Bacterial distribution by body site. This figure shows the distribution by body site of bacteria that have been sequenced under the HMP or are in the sequencing pipelines.



# Global patterns of 16S rRNA diversity at a depth of millions of sequences per sample

J. Gregory Caporaso<sup>a</sup>, Christian L. Lauber<sup>b</sup>, William A. Walters<sup>c</sup>, Donna Berg-Lyons<sup>b</sup>, Catherine A. Lozupone<sup>a</sup>, Peter J. Turnbaugh<sup>d</sup>, Noah Fierer<sup>b,e</sup>, and Rob Knight<sup>a,f,1</sup>

Target gene:



Amplification primers with annealing sites:



# MSA after forward primer

Jalview 2.7  
File Tools Vamsas Help Window

C:\Users\ranjit\Desktop\morrow-working-files\RDP\bacterial16S\_508\_mod5.stk

File Edit Select View Format Colour Calculate Web Service

750 760 770 780 790 800 810 820 830 840 850 860 870 880

NC\_007292/1-1566 C G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G G U G C G A G C G U U A U C G G A A U U A C U G G G C G U A A A . G A G U A C G U A G G U G G U . U U G U U A A G U C A G . A U G U G . A A A U C C C G U A G C U C A A C U U A G G A . A C U G C A U U U G  
NC\_008769/1-1532 C G U G C C A G C A G C C G C G G U A A U A C G U A G . . . . . G G U G C G A G C G U U G U C C G G A A U U A C U G G G C G U A A A . G A G C U C G U A G G U G G U . U U G U C G C G U U G U . U C G U G . A A A U C U C A G C G G C U U A A C U G U G A G . C G U G C G G G C G  
NC\_008800/1-1543 C G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G G U G C A A G C G U U A A U C G G A A U U A C U G G G C G U A A A . G C G C A C G C A G G C G G U . U U G U U A A G U C A G . A U G U G . A A A U C C C G C G C U U A A C G U G G G A . A C U G C A U U U G  
NC\_009446/1-1533 C G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G G U G C A A G C G U U A U C C G G A A U G A C U G G G C G U A A A . G C G C A C G C A G G U G G U . U U U A U A A G U C A G . G U G U G . A A A U C C C U G G G C U C A A C C U A G G A . A U U G C A U U U G  
NC\_008767/1-1541 C G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G G U G C G A G C G U U A U C C G G A A U U A C U G G G C G U A A A . G C G G G C G C A G A C G G U . U A C U U A A A G C A G G . A U G U G . A A A U C C C G G G C U C A A C C C G G G A . A C U G C G U U C U  
NC\_009445/1-1489 C G U G C C A G C A G C C G C G G U A A U A C G A A G . . . . . G G G G C U A G C G U U U C C G G A A U A C U G G G C G U A A A . G G G U G C G U A G G C G G G . U C U U U A A G U C A G . G G G U G . A A A U C C U G G A G C U A A C C U C A G A . A C U G C C U U U G  
NC\_009443/1-1549 C G U G C C A G C A G C C G C G G U A A U A C G U A G . . . . . G U C C G A G C G U U G U C C G G A U U U A U G G G C G U A A A . G C G A G C G C A G G C G G U . U U G A U A A G U C U G . A A G U A . A A A G G C U G U G G C U U A A C C A U A G U . A C . C C U U U G G  
NC\_009442/1-1549 C G U G C C A G C A G C C G C G G U A A U A C G U A G . . . . . G U C C G A G C G U U G U C C G G A U U U A U G G G C G U A A A . G C G A G C G C A G G C G G U . U U G A U A A G U C U G . A A G U A . A A A G G C U G U G G C U U A A C C A U A G U . A C . G C U U U G G  
NC\_009441/1-1514 C G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G A U C C A A G C G U U A U C C G G A A U C A U U G G G U U U A A A . G G G C A G C G U A G G C G G U . U U A G U A A G U C A G . U G G U G . A A A G C C A U C G C U C A A C C U G G G A . A C G C C A U U G  
NC\_009049/1-1467 C G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G G G G C U A G C G U U A U C C G G A A U U A C U G G G C G U A A A . G C G C A C G U A G G C G G U . U C G G A A A G U C A G . A G G U G . A A A U C C A G G G C U C A A C C C U G G A . A C U G C C U U U G  
NC\_003454/1-1520 C G U G C C A G C A G C C G C G G U A A U A C G U A U . . . . . G U C A C G A G C G U U A U C C G G A U U U A U G G G C G U A A A . G C G C U C U A G G G C G G U . U A U G U A A G U C U G . A U G U G . A A A U G C A G G G C U C A A C C U C U G U A . - U U G C U U G G  
NC\_008369/1-1528 C G U G C C A G C A G C C G C G G U A A U A C G G G G . . . . . G G U G C A A G C G U U A A U C G G A A U U A C U G G G C G U A A A . G G G U C G U A G G U G G U . U U G U U A A G U C A G . A U G U G . A A A G C C C A G G G C U C A A C C U U G G A . A C U G C A U U U G  
NC\_007722/1-1486 C G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G G A G C U A G C G U U G U U C C G A A U U A C U G G G C G U A A A . G C G C G C U A G G C G G C . U A U U U A A G U C A G . G G G U G . A A A U C C C G G G G C U C A A C C C C G G A . A C U G C C U U U G  
NC\_008009/1-1502 U G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G G U G C A A G C G U U A U C C G G A A U U A U G G G C G U A A A . G G G C G U A G G C G G U . U A G U A A A G U C U C . U A G U G . A A A U C C C G G G C U C A A C C U C G G A . C C U G C U A G G G  
NC\_003450/1-1524 C G U G C C A G C A G C C G C G G U A A U A C G U A G . . . . . G G U G C G A G C G U U G U C C G G A A U U A C U G G G C G U A A A . G A G C U C G U A G G U G G U . U U G U C G C G U C G U . C U G U G . A A A U C C C G G G C U U A A C U U C G G G . C G U G C A G G C G  
NC\_002771/1-1525 U G U G C C A G C A G C C G C G G U A A U A C A U A G . . . . . G G U G C A A G C G U U A U C C G A A A U U A U G G G U G U A A . G A G U U C G U A G G U U G U . U U G U U A A G U C A G . A A G U U . A A A U C C C G G G C U C A A C C C U G G C . C C . G C U U U U G  
NC\_005966/1-1508 U G U G C C A G C A G C C G C G G U A A U A C A G A G . . . . . G G U G C A A G C G U U A A U C G G A A U U A C U G G G C G U A A A . G C G C G C U A G G C G G C . C A A U U A A G U C A A . A U G U G . A A A U C C C G A G C U U A A C U U G G G A . A U U G C A U U C G  
NC\_009439/1-1536 C G U G C C A G C A G C C G C G G U A A U A C G A A G . . . . . G G U G C A A G C G U U A A U C G G A A U U A C U G G G C G U A A A . G C G C G C U A G G U G G U . U C G U U A A G U U G G . A U G U G . A A A G C C C G G G C U C A A C C U G G G A . A C U G C A U C C A  
NC\_009438/1-1543 C G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G G U C C A A G C G U U A U C C G A A U U A C U G G G C G U A A A . G C G U G C G A G G C G G U . U U G U U A A G C A G . A U G U G . A A A C C U G G G C U C A A C C U G G A . A U A G C A U U U C  
NC\_009437/1-1544 C G U G C C A G C A G C C G C G G U A A U A C G U A G . . . . . U G U G C G A G C G U U G U C C G G A A U U A C U G G G C G U A A A . G G G U C G U A G G C G G C . U A U G C G A G U U A A . G C G U G . A A A G C C U U A G G G C U C A A C C U A A G G . A U U G C G C U U A  
NC\_004129/1-1539 U G U G C C A G C A G C C G C G G U A A U A C A G A G . . . . . G G U G C A A G C G U U A A U C G G A A U U A C U G G G C G U A A A . G C G C G C U A G G U G G U . U U G U U A A G U U G G . A U G U G . A A A G C C C G G G G C U C A A C C U G G G A . A C U G C A U C C A  
NC\_009436/1-1540 C G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G G U G C A A G C G U U A A U C G G A A U U A C U G G G C G U A A A . G C G C A G C G U A G G C G G U . C U G U C A A G U C G G . A U G U G . A A A U C C C G G G G C U C A A C C U G G G A . A C U G C A U U C G  
NC\_009434/1-1537 C G U G C C A G C A G C C G C G G U A A U A C G A A G . . . . . G G U G C A A G C G U U A A U C G G A A U U A C U G G G C G U A A A . G C G C G C U A G G U G G U . U C G U U A A G U U G G . A U G U G . A A A G C C C G G G C U C A A C C U G G G A . A C U G C A U C C A  
NC\_008268/1-1518 C G U G C C A G C A G C C G C G G U A A U A C G U A G . . . . . G G U G C A A G C G U U A U C C G G A A U U A C U G G G C G U A A A . G A G U U C G U A G G C G G U . U U G U C G C U C G U . U U G U G . A A A A C U C A C A G C U C A A C U G U G A G . C C U G C A G G C G  
NC\_008752/1-1529 C G U G C C A G C A G C C G C G G U A A U A C G U A G . . . . . G G U G C A A G C G U U A A U C G G A A U U A C U G G G C G U A A A . G C G U G C G C A G G C G G U . G A U G U A A G A C A G . A U G U G . A A A U C C C G G G G C U C A A C C U G G G A . A C U G C A U U U G  
NC\_008751/1-1549 C G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G G U G C G A G C G U U A A U C G G A A U A C U G G G C G U A A A . G C G C A C G U A G G C U G C . U U G G U A A G U C A G . G G G U G . A A A G C C C G G G C U C A A C C C G G G A . A U U G C C U U U G  
NC\_007712/1-1542 C G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G G U G C A A G C G U U A A U C G G A A U U A C U G G G C G U A A A . G C G U C A A G C G G U . C U G U U A A G U C A G . A U G U G . A A A U C C C G G G G C U U A A C C U G G G A . A C U G C A U U U G  
NC\_008750/1-1543 C G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G G U C C A A G C G U U A A U C G G A A U U A C U G G G C G U A A A . G C G U G C G A G G C G G U . U U G U U A A G C A G . A U G U G . A A A G C C C G G G C U C A A C C U A G G A . A U A G C A U U U C  
NC\_008359/1-1455 C G U G C C A G C A G C C G C G G U A A U A C G A A G . . . . . G G G G C U A G C G U U G U C C G G A A U U A C U G G G C G U A A A . G C G C A C G U A G G C G G A . C U U U U A A G U C A G . G U G U G . A A A U C C G A G G G C U C A A C C U C G G A . A C U G C A U U U G  
NC\_004088/1-1543 C G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G G U G C A A G C G U U A A U C G G A A U U A C U G G G C G U A A A . G C G C A C G C A G G C G G U . U U G U U A A G U C A G . A U G U G . A A A U C C C G C G C U U A A C G U G G G A . A C U G C A U U U G  
NC\_007677/1-1540 C G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G G U G C A A G C G U U G U C C G G A A U A C U G G G U G U A A A . G G G U G U G C A G G C G G G . G C A G C A A G U C G G . A U G U G . A A A C C C A U G G G C U U A A C C A U G G A . G G U G C A U U C G  
NC\_003047/1-1485 C G U G C C A G C A G C C G C G G U A A U A C G A A G . . . . . G G G G C U A G C G U U G U U C C G G A A U U A C U G G G C G U A A A . G C G C A C G U A G G C G G A . U U G U U A A G U G A G . G G G U G . A A A U C C A G G G C U C A A C C C U G G A . A C U G C C U U U C  
NC\_005957/1-1552 C G U G C C A G C A G C C G C G G U A A U A C G U A G . . . . . G U G G C A A G C G U U A U C C G G A A U U A U U G G G C G U A A A . G C G C G C G A G G U G G U . U U C U U A A G U C U G . A U G U G . A A A G C C A C G G C U C A A C C U G G A . G G G U C A U U G G  
NC\_005956/1-1488 C G U G C C A G C A G C C G C G G U A A U A C G A A G . . . . . G G G G C U A G C G U U G U U C C G G A A U U A C U G G G C G U A A A . G C G C A U G U A G G C G G A . U A U U U A A G U C A G . A G G U G . A A A U C C A G G G C U C A A C C C U G G A . A C U G C C U U U G

secondary structure  
reference positions C G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G g g g C a a g C g u u g u c C G G A A U U A c U G G G C G U A A A . G a G c G c G a G g G g G c . c g c c a A g u c g g . g u G c g . A A A u u c c g g g G c U u A A C c c g g g a . A a c G C a c c c g

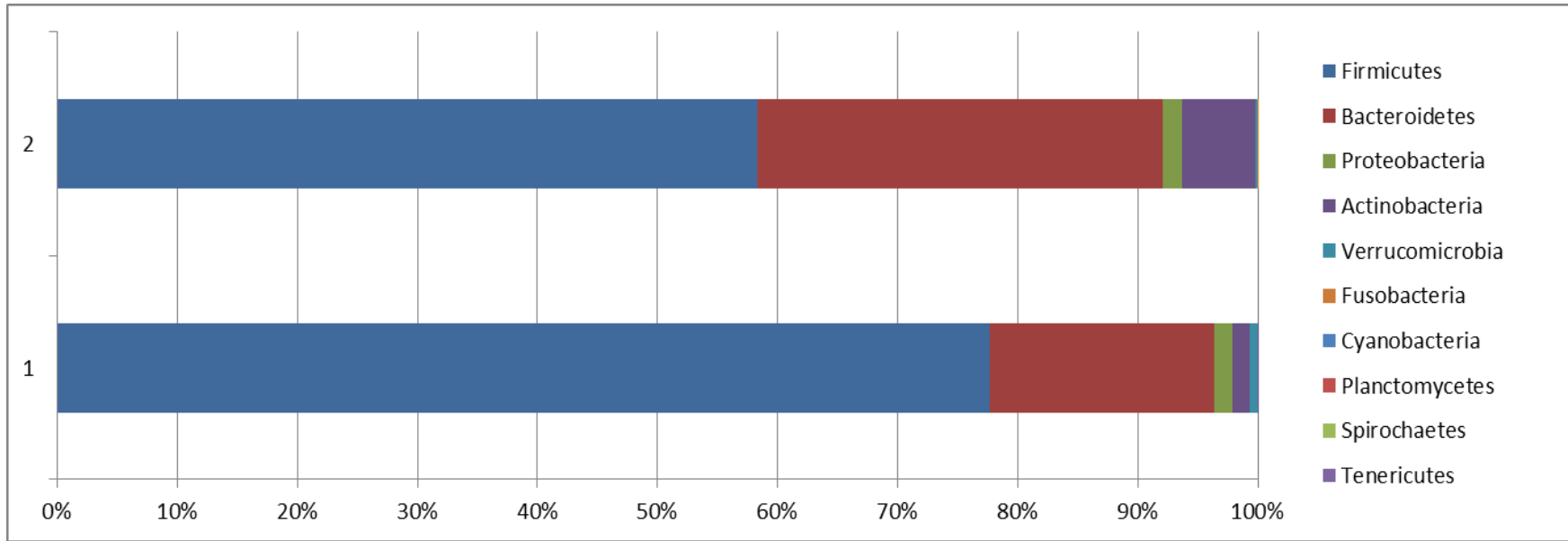
Consensus C G U G C C A G C A G C C G C G G U A A U A C G G A G . . . . . G G U G C A A G C G U U A U C G G A A U U A C U G G G C G U A A A . G C G C G C U A G G C G G U . U U G U U A A G U C A G U A U G U G . A A A U C C C G G G C U C A A C C U G G G A . A C U G C A U U U G

Sequence 13 ID: NC\_007722 Nucleotide: Uracil (561)

1:50 PM 10/10/2011

# Microbiome at UAB

Normal Diabetic



OPEN ACCESS Freely available online

PLoS one

## Gut Microbiota in Human Adults with Type 2 Diabetes Differs from Non-Diabetic Adults

Nadja Larsen<sup>1\*</sup>, Finn K. Vogensen<sup>1</sup>, Frans W. J. van den Berg<sup>1</sup>, Dennis Sandris Nielsen<sup>1</sup>, Anne Sofie Andreassen<sup>2</sup>, Bente K. Pedersen<sup>2</sup>, Waleed Abu Al-Soud<sup>3</sup>, Søren J. Sørensen<sup>3</sup>, Lars H. Hansen<sup>3</sup>, Mogens Jakobsen<sup>1</sup>

<sup>1</sup> Department of Food Science, University of Copenhagen, Frederiksberg, Denmark, <sup>2</sup> Department of Infectious Diseases and CMRC, University Hospital Rigshospitalet, Copenhagen, Denmark, <sup>3</sup> Department of Biology, University of Copenhagen, Copenhagen, Denmark

The proportions of phylum Firmicutes and class Clostridia were significantly reduced in the diabetic group compared to the control group (P = 0.03).

# Sequencing RNA

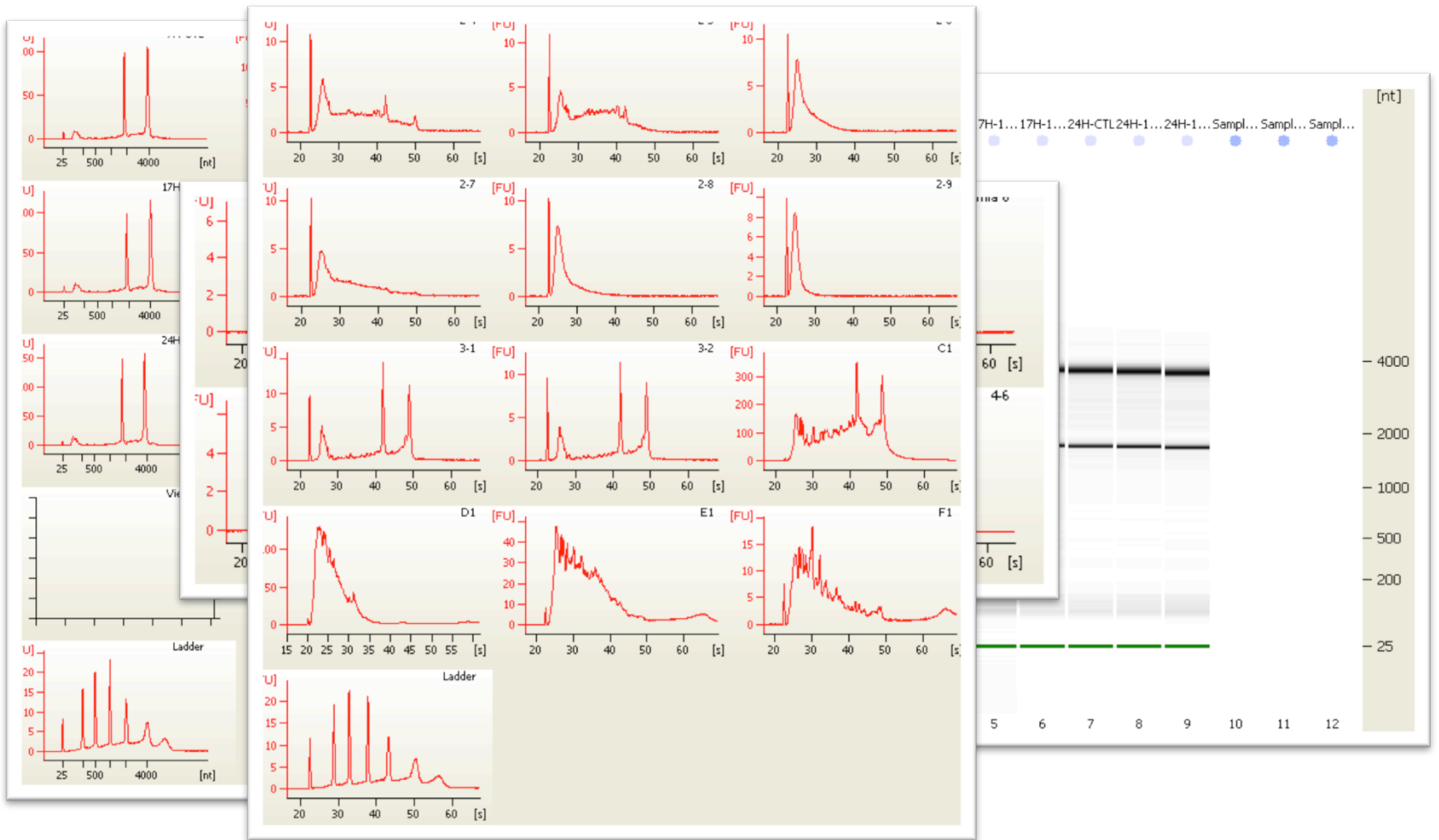
# RNA Applications

- mRNA Sequencing (RefSeq, RNASeq)
- microRNA Sequencing
- RNA-IP-Sequencing
- CLIP or HITS-CLIP or PAR-CLIP
- Ribosome Profiling

# Advantages of RNA-Seq

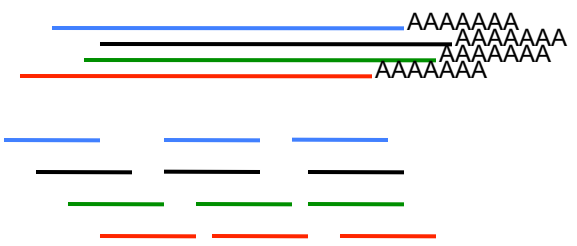
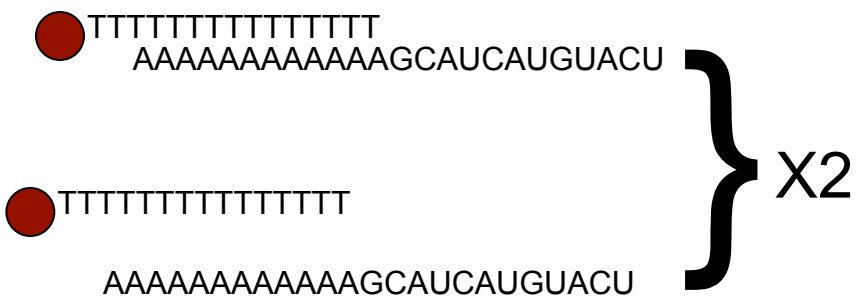
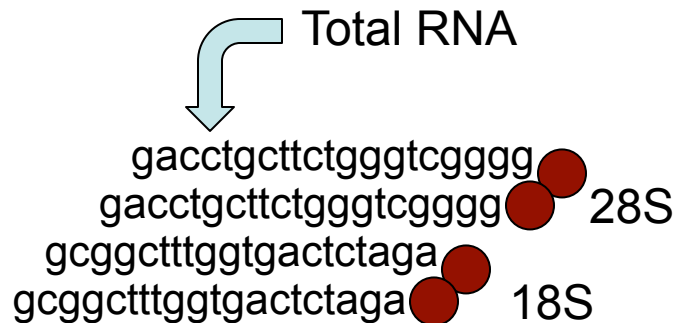
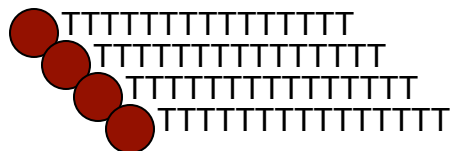
- Digital gene expression
  - Simply count the number of reads for a given transcript
- Greater dynamic range
- No hybridization bias
- Not dependent on known content
- Generate alternative splice/exon usage
- Identify variants
- Allele Specific Expression
- Identify RNA editing

# RNA Quality



# RNA-Seq Library Prep

Total RNA  Mix/65°C/25°C



Fragment

↓

Random Priming to make cDNA

↓

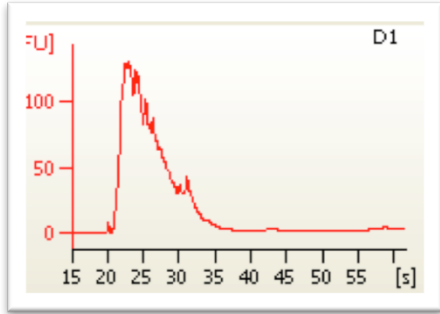
Random Priming to make cDNA



Sequence ready library



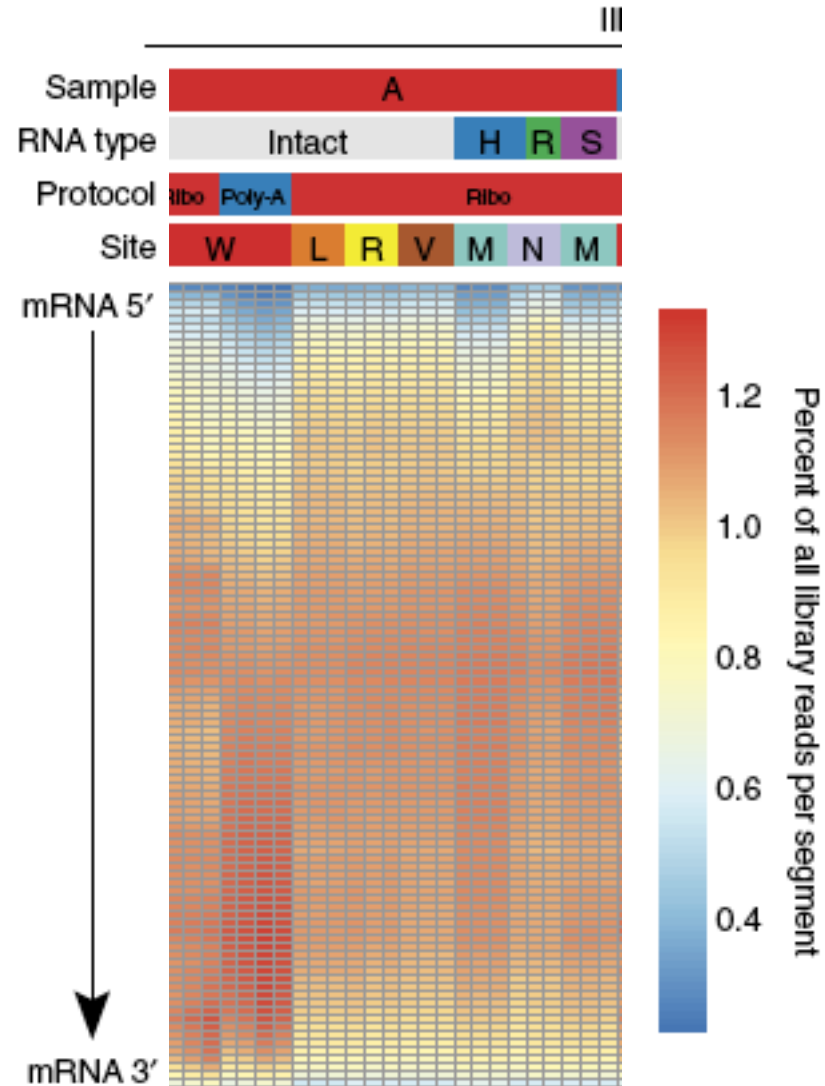
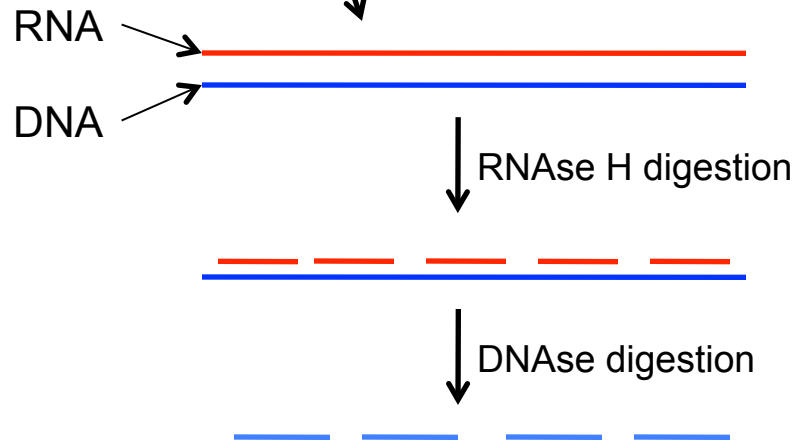
# Degraded RNA (FFPE)



Total RNA

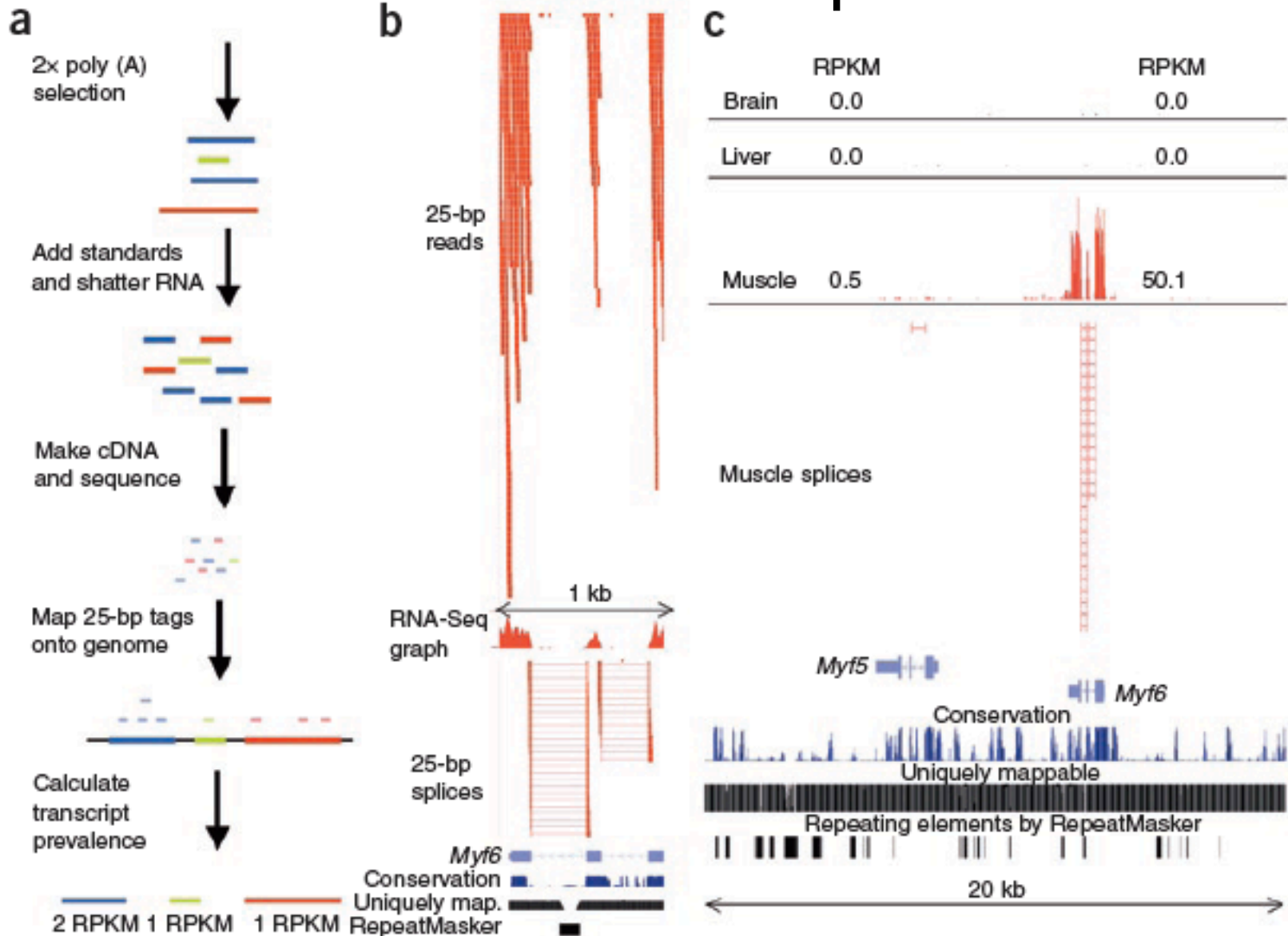
gacctgcttctgggtcgggg  
gacctgcttctgggtcgggg 28S

gcggctttggtgactctaga  
gcggctttggtgactctaga 18S



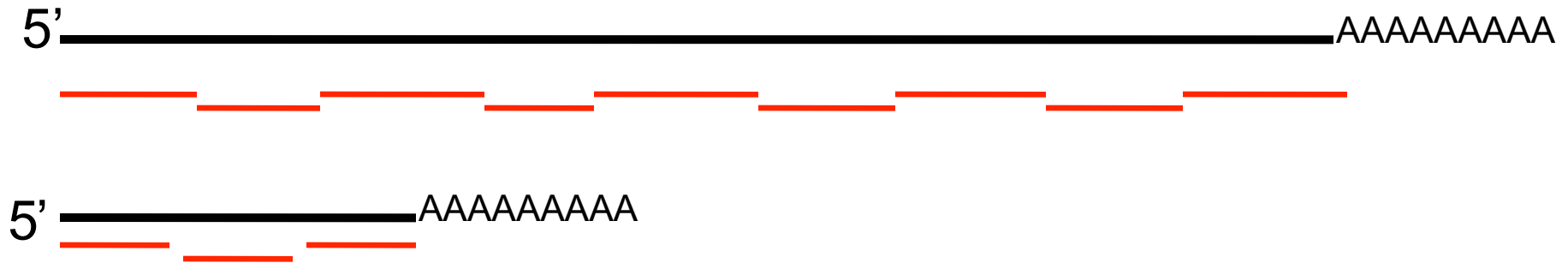
From Li et al., 2014 Nat. Biotech.32:915

# mRNA-Seq



# Digital Gene Expression

- Caveat?

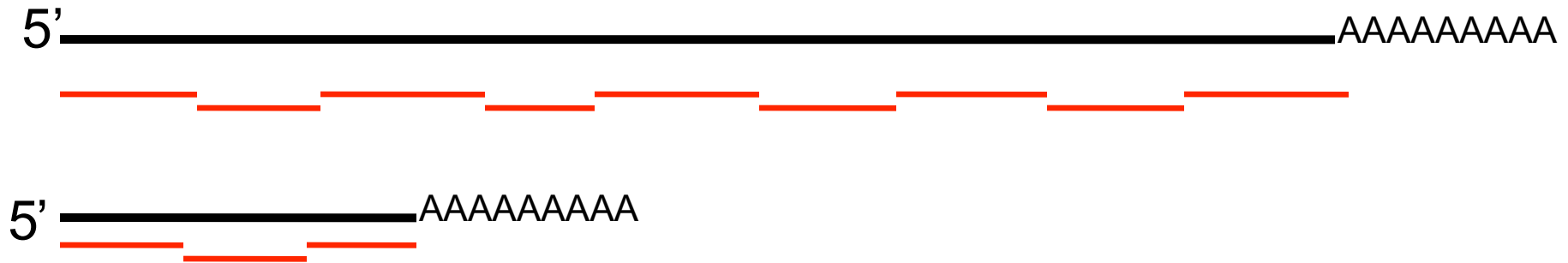


Gene 1 has 9 reads  
Gene 2 has 3 reads

Gene 1 would appear to be expressed at  
3X the amount of Gene 2

# Digital Gene Expression

- Caveat?



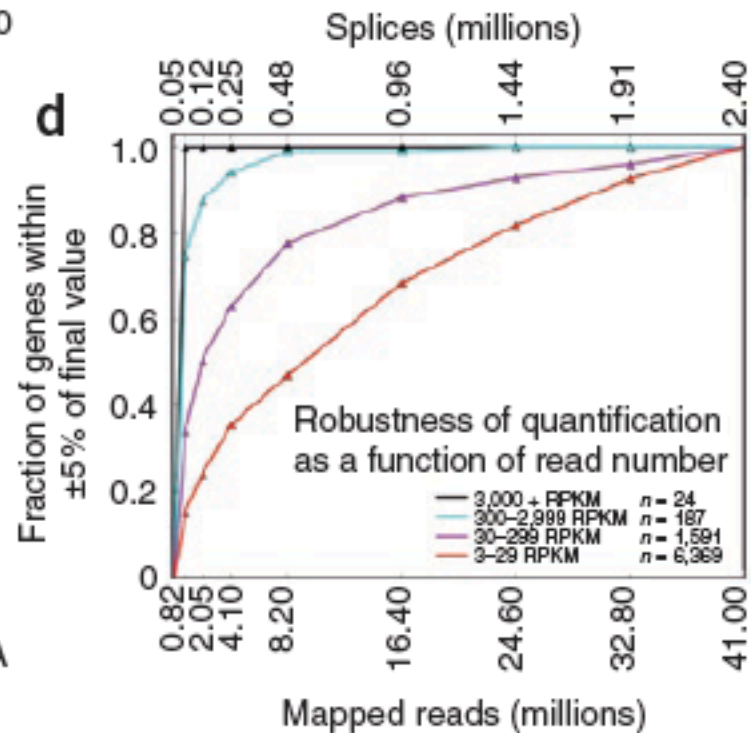
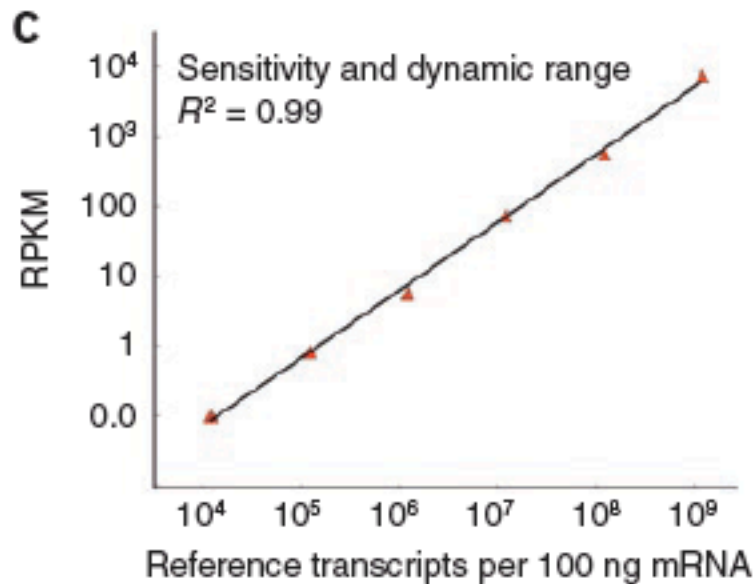
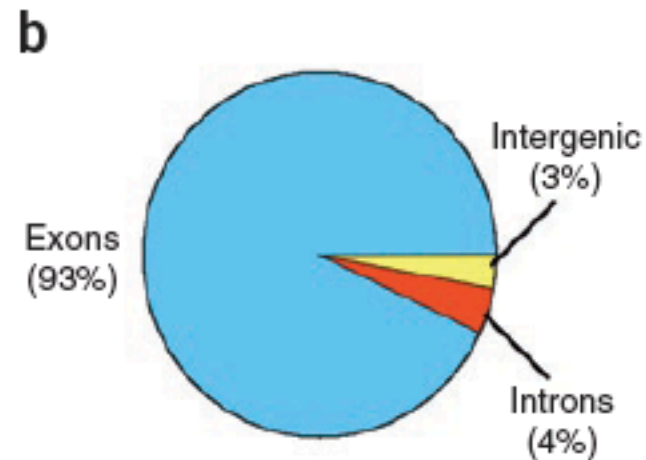
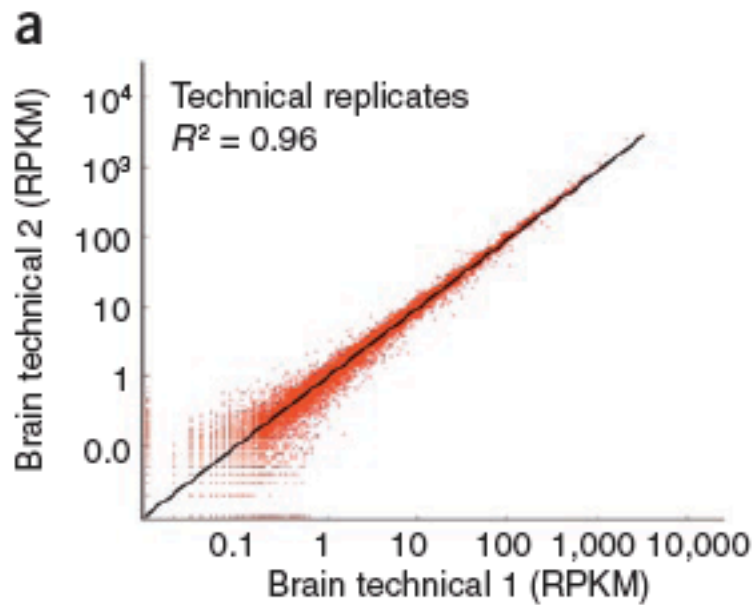
Gene 1 has 9 reads  
Gene 2 has 3 reads

Gene 1 is 9kb  
Gene 2 is 3kb

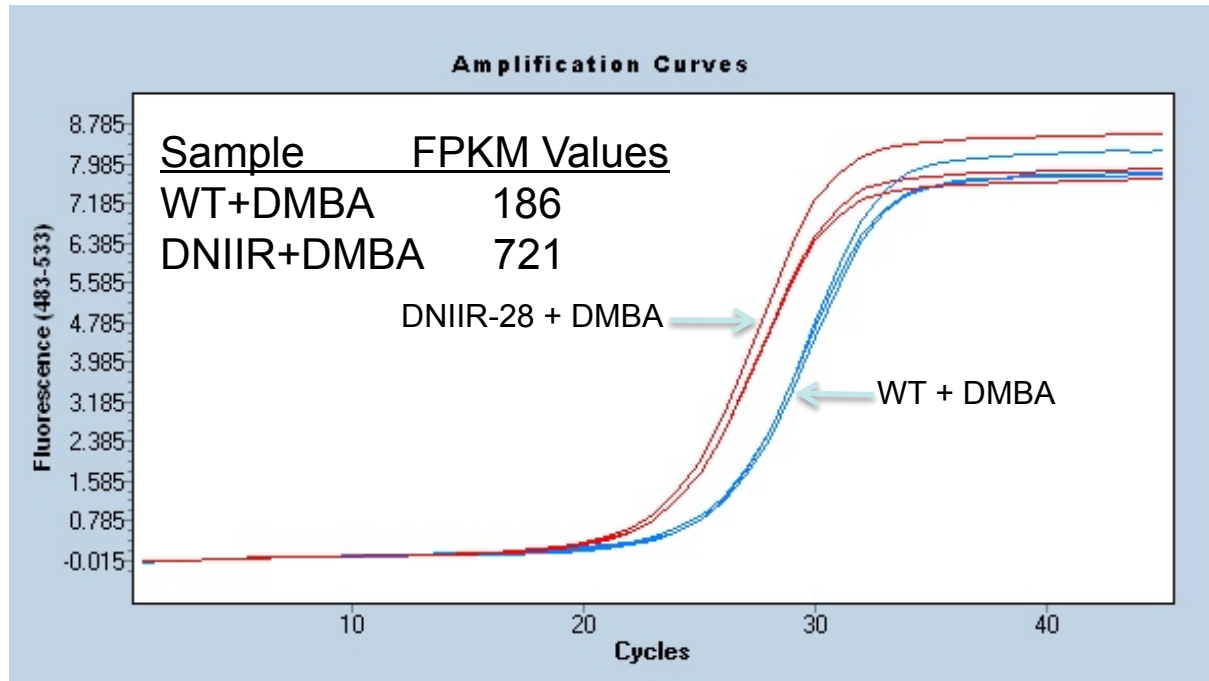
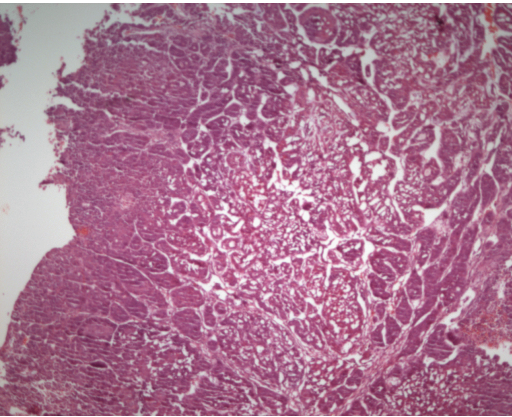
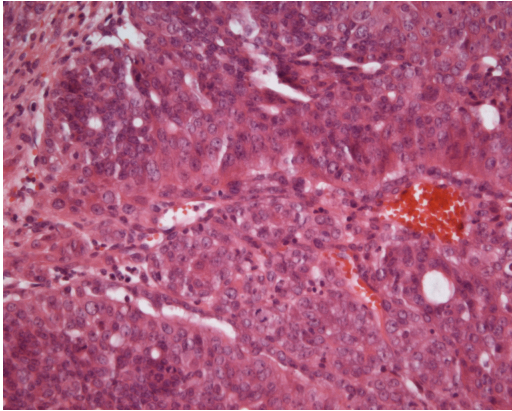
Normalize by length  
e.g.  $9\text{reads}/9\text{kb}=1$

RPKM: reads per kilobase of exon model per million mapped reads

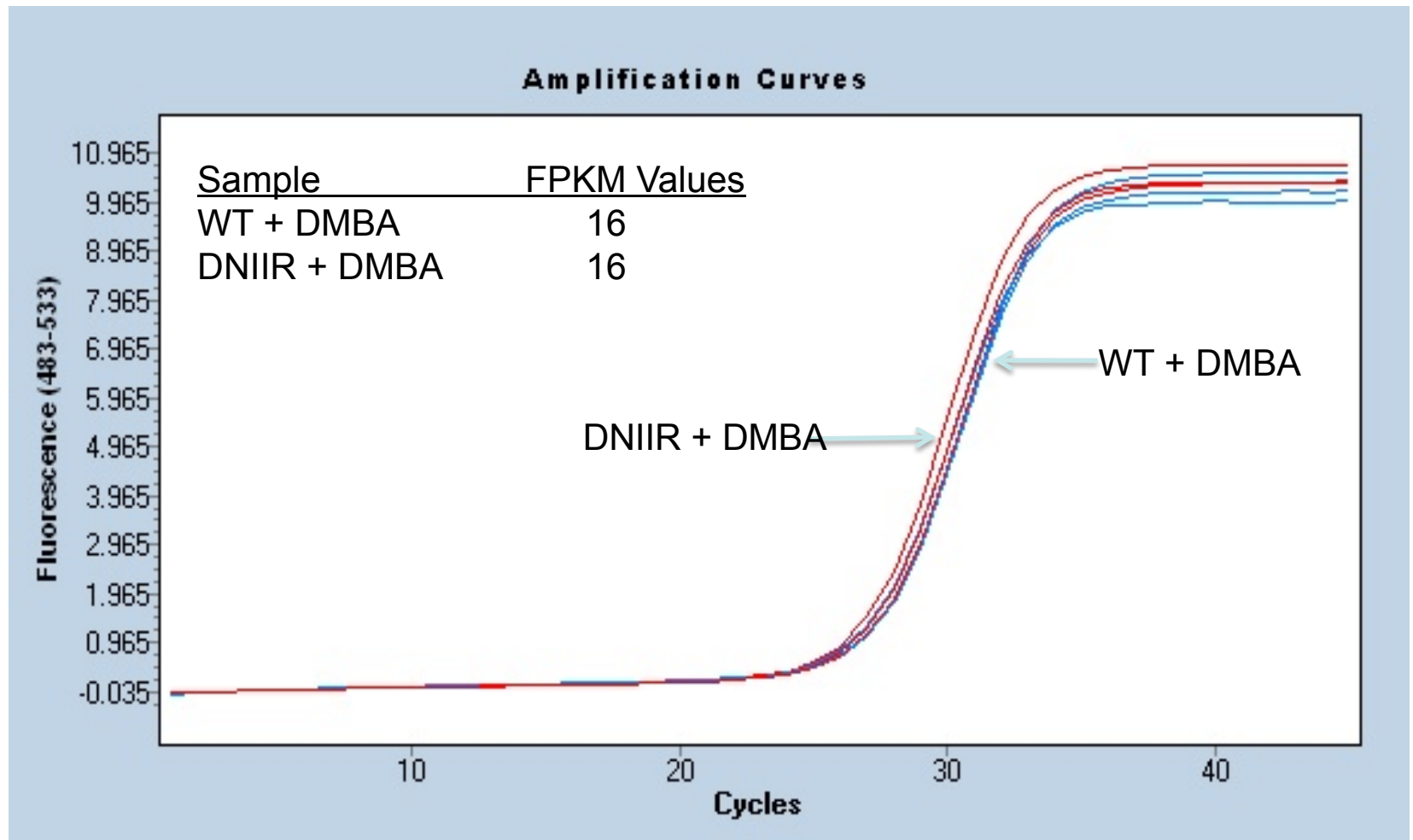
FPKM: fragment of reads per kilobase of exon model per million mapped reads (usually 25bp fragments).



# Keratin 8

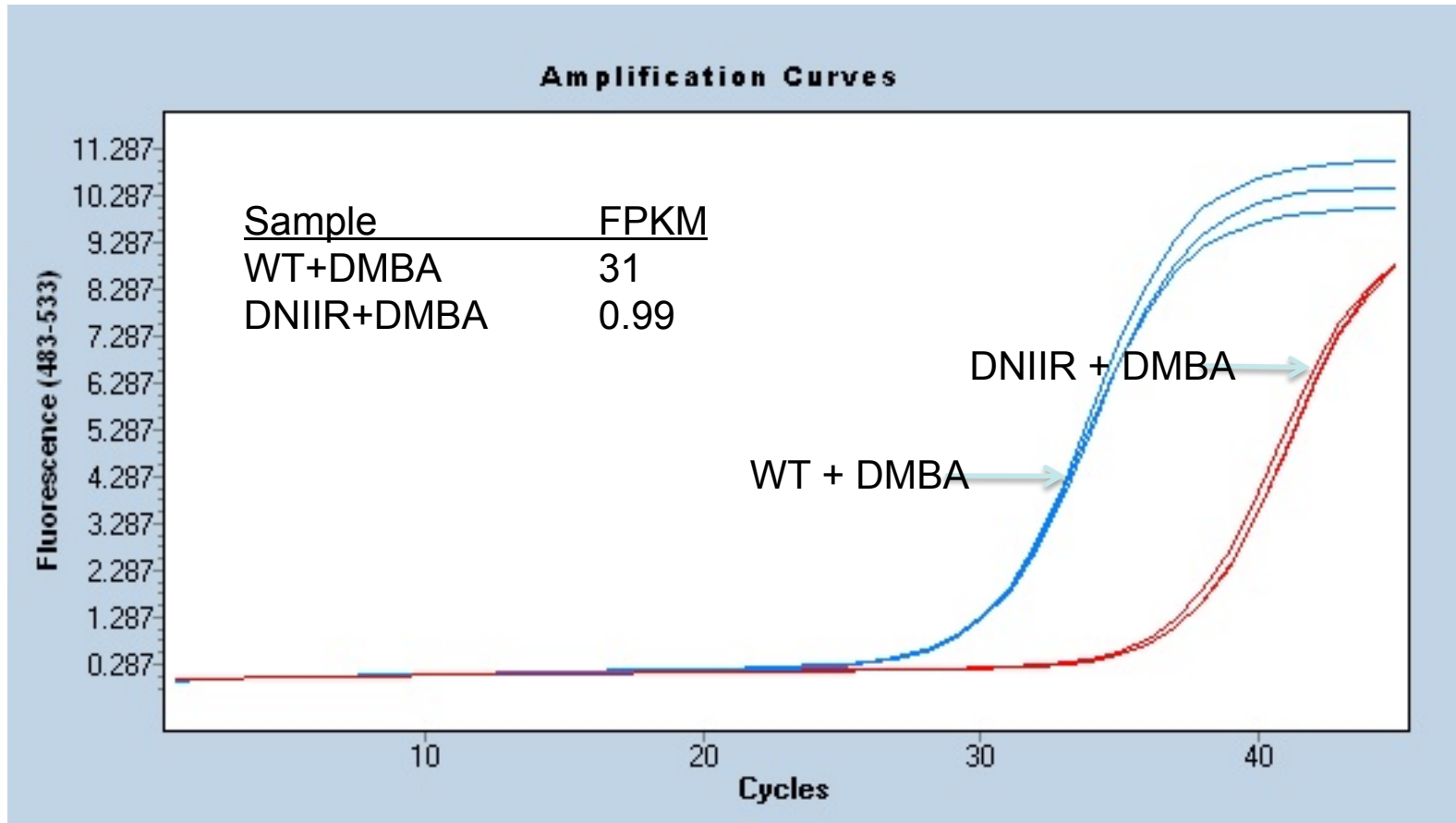


# Lipase Maturation Factor 1 (Lmf1)

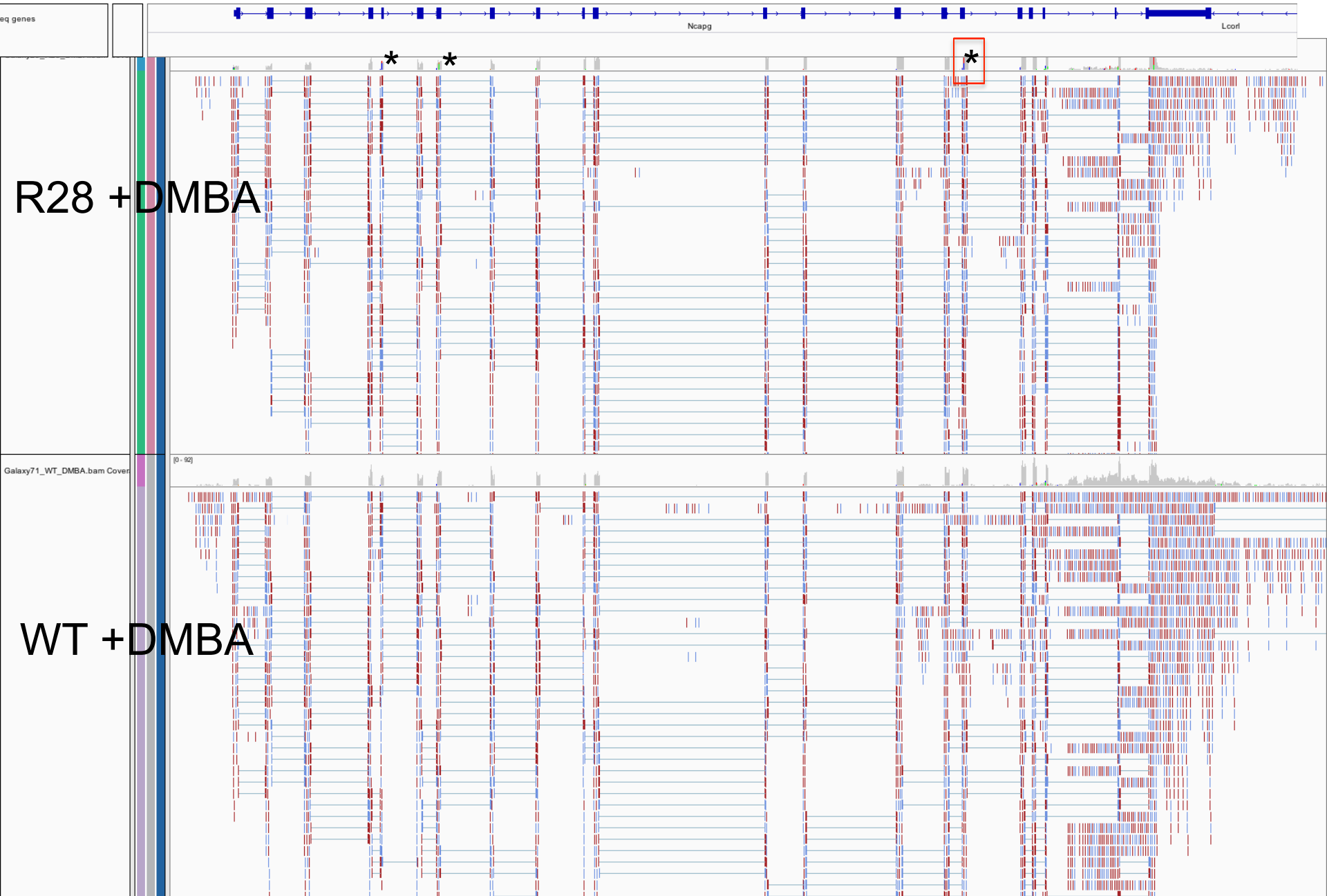




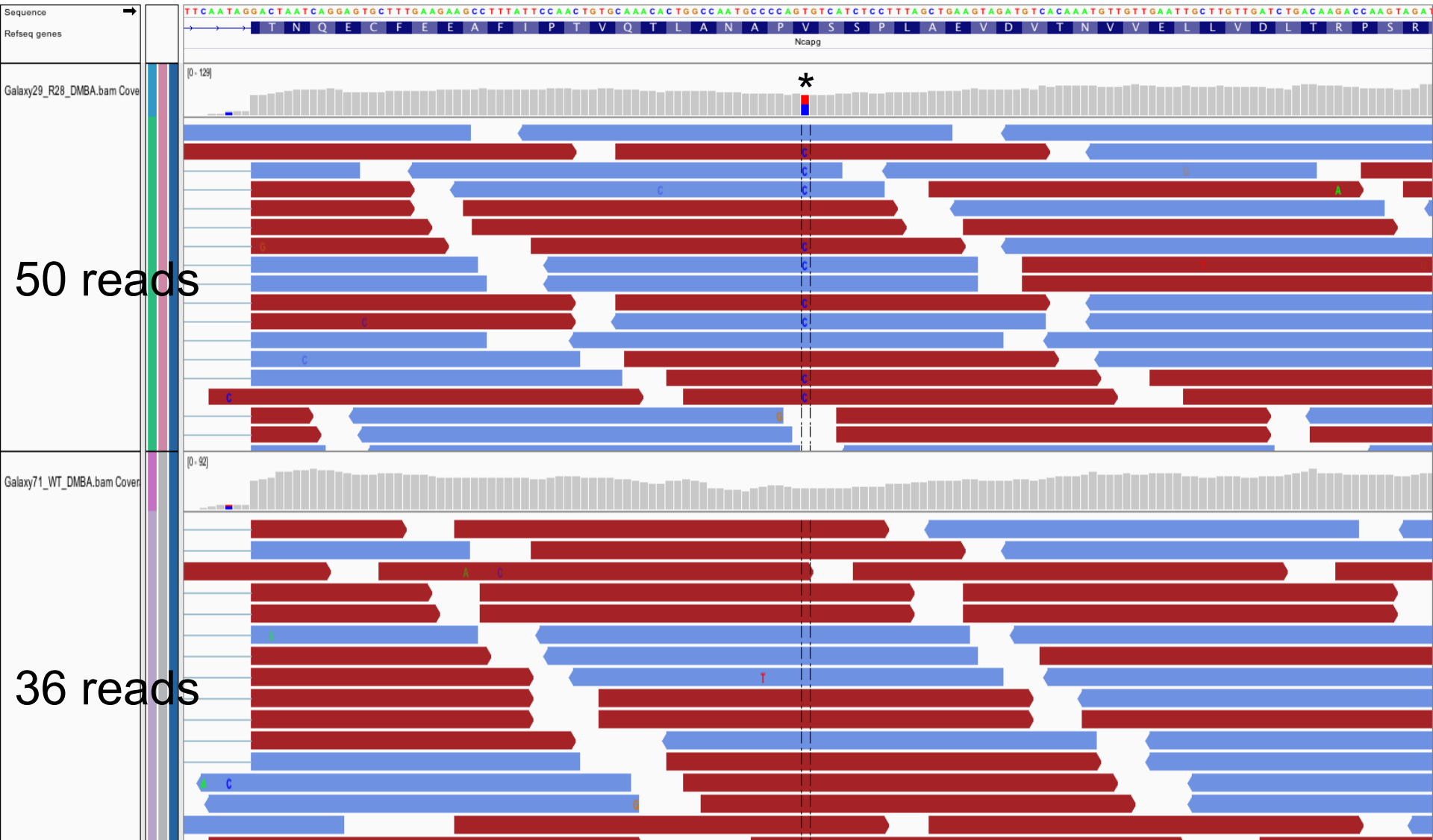
# Lysophosphatidic acid receptor 3



# Ncapg: Non-SMC condensin I complex, subunit G

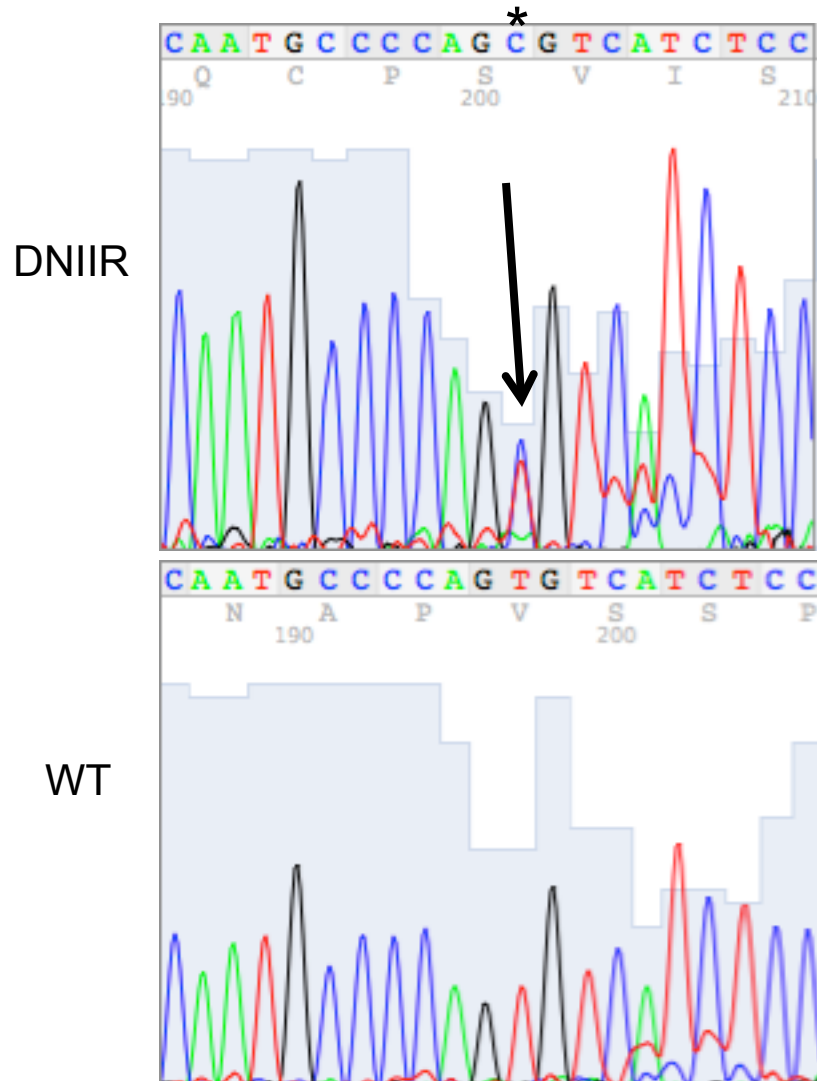


# Exon 16 of Ncapg



T-C mutation resulting in a Val-Ala change in the protein

# Sequence Confirmation of Ncapg mutation

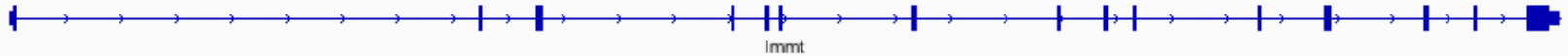


T>C mutation resulting in an Ala>Val change at position aa784 in the protein. The other mutations were a polymorphic T>C change at aa242 and an A>G change at aa347 resulting in a non-synonymous change from Arg>Lys.

# Alternative Exon Usage



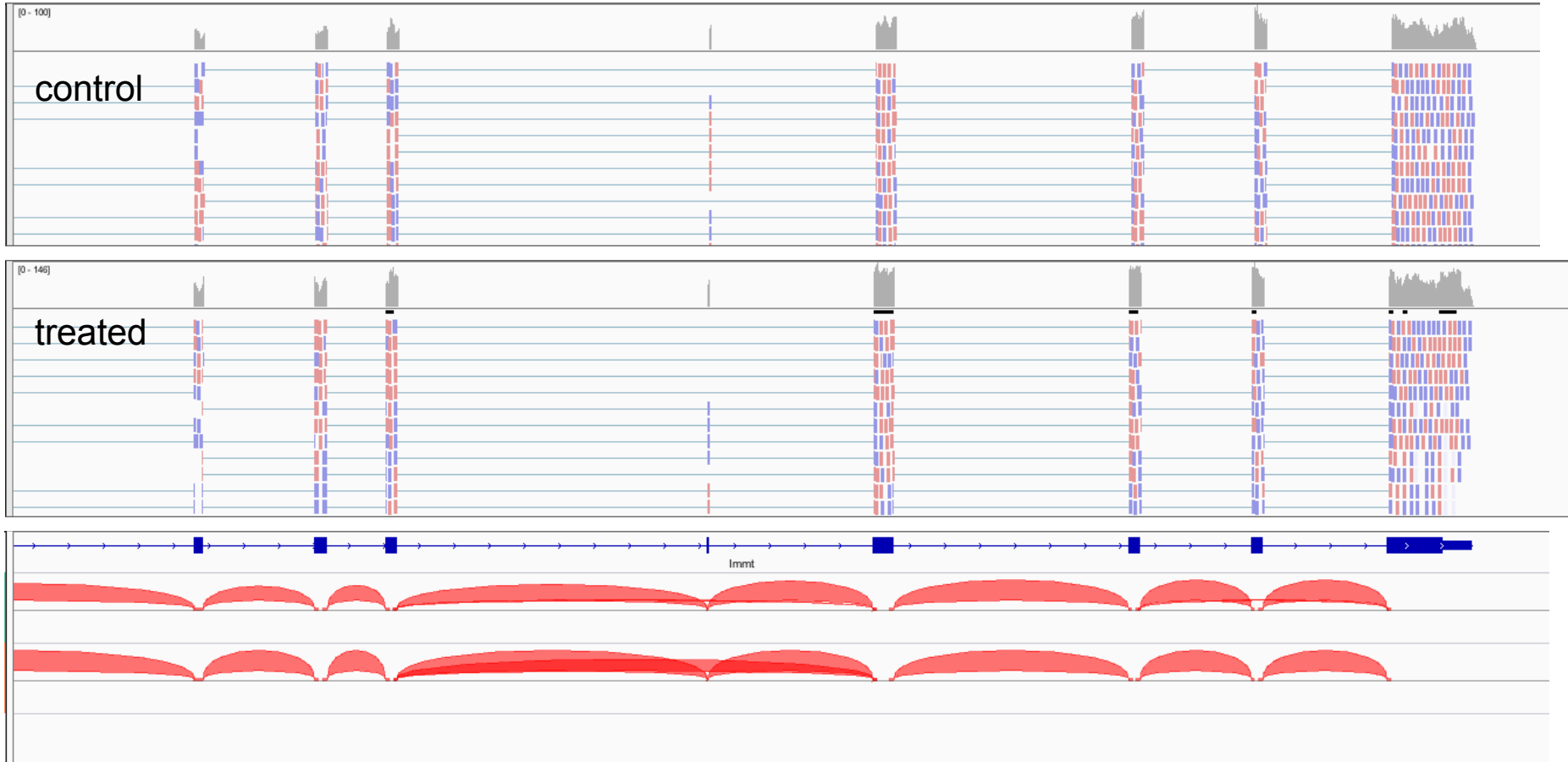
Trim24



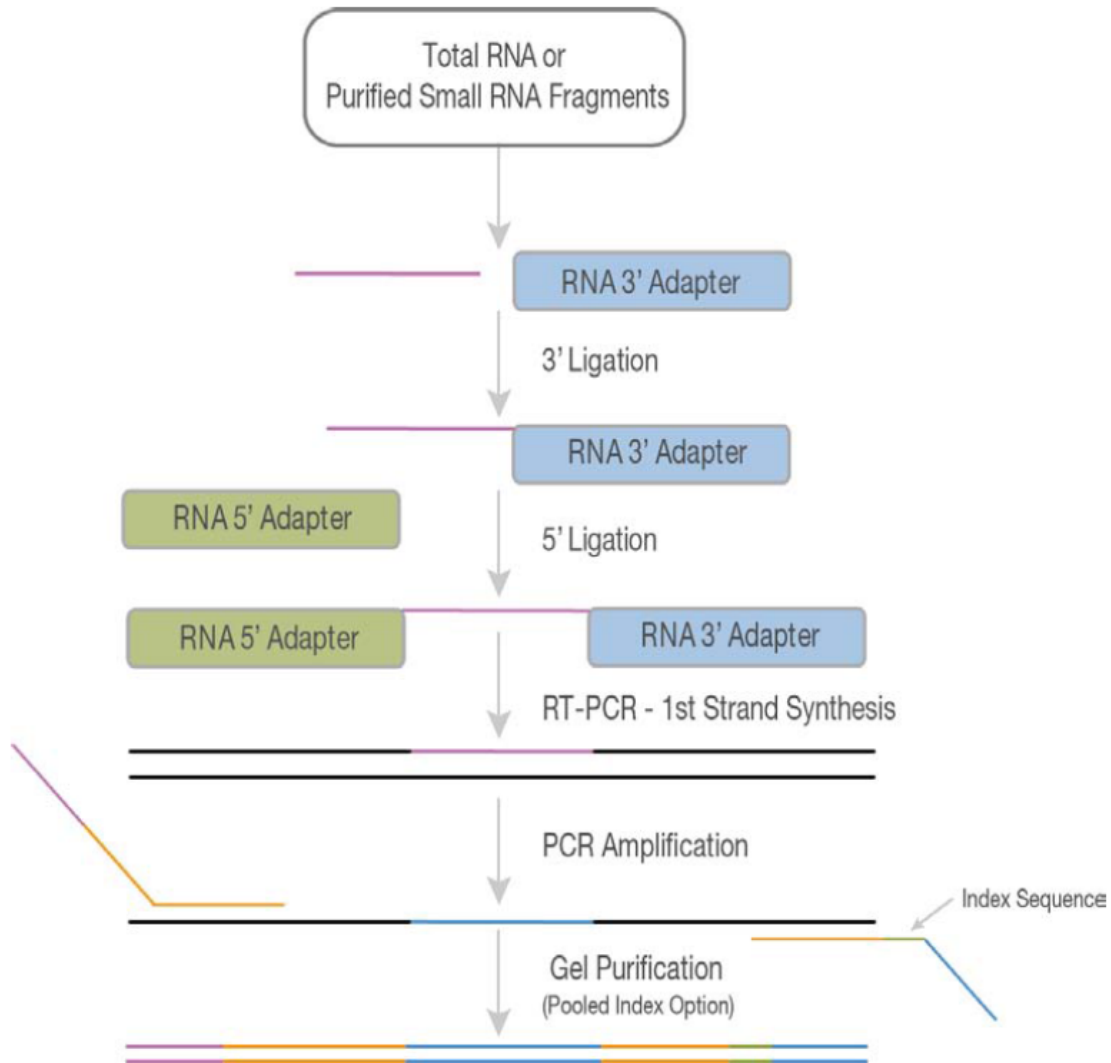
Immt



# Alternative splicing

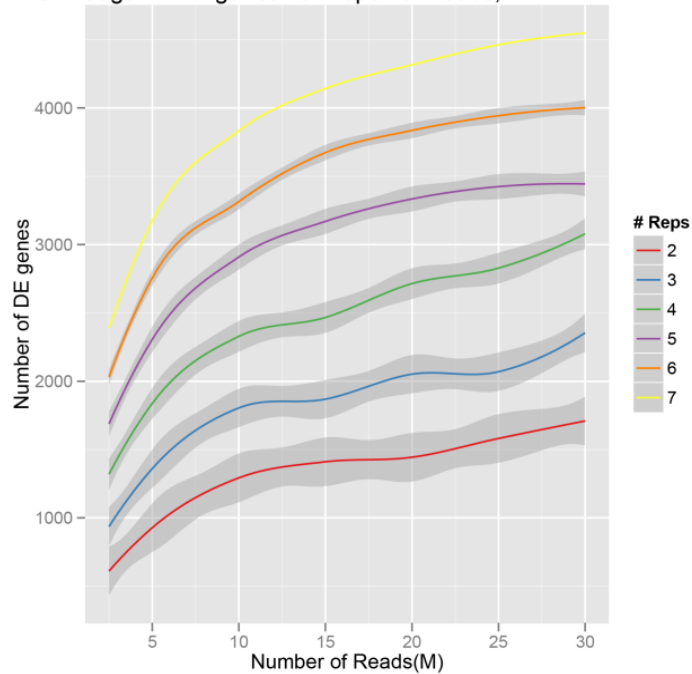


# microRNA Seq

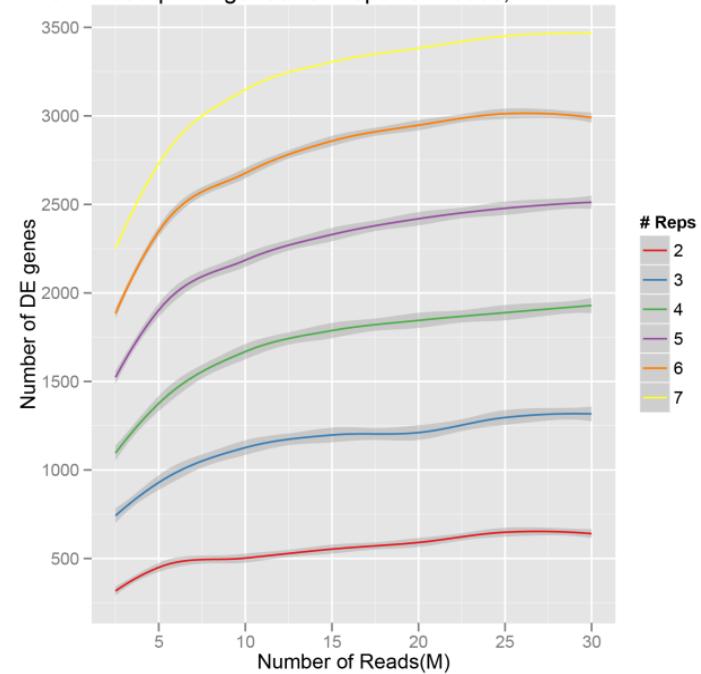


# Sequencing Depth v. Biological Replicates

**a** edgeR #DE genes vs. Repls vs. Reads, FDR 0.01

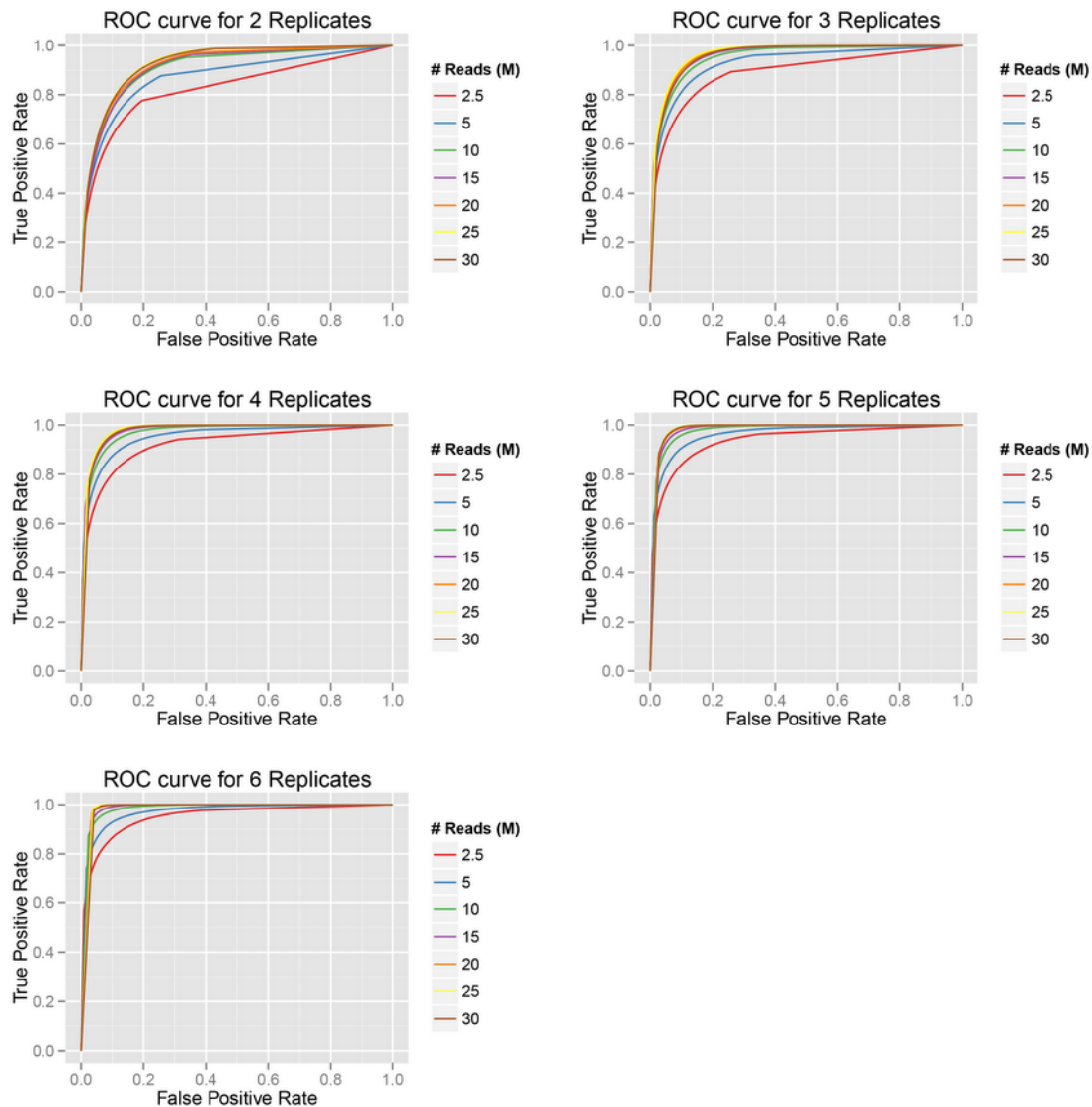


**b** DESeq #DE genes vs. Repls vs. Reads, FDR 0.05



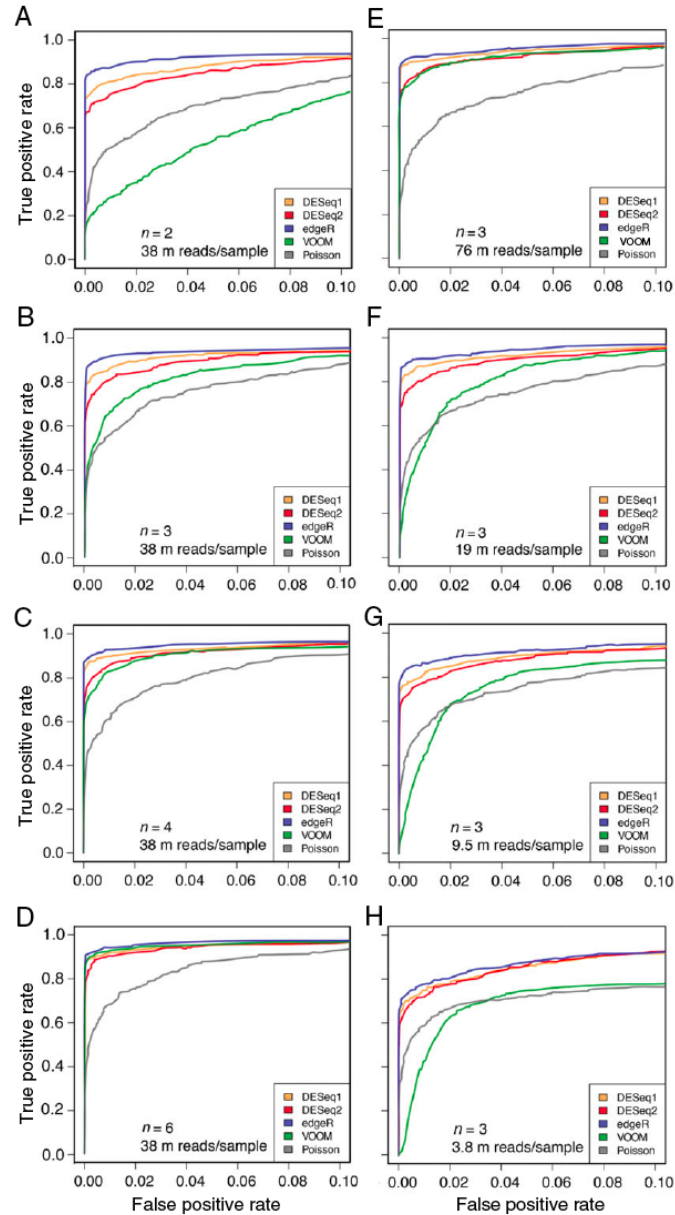


# ROC curves with Replicates v. Depth



# Replicates v. Depth with different DE packages

Williams AG, Thomas S., Wyman SL., Holloway AH. 2014.  
Curr Prot Human Genet 11.13.1



# Recommendations

1. Biological replication is important
  - The more you can afford the better your data
2. Increasing read depth does not substitute for increasing replicates
3. The type of experiment matters
  - Cell lines or inbred strains vs. outbred populations (humans for example).

# Summary

- Several different platforms exist utilizing different technologies.
- Generate between 500 million to 600 Billion bases of sequence information per run.
- Several applications including Whole genome sequencing, Targeted genomic seq., ChIP-Seq and mRNA-Seq, among others.
- Data files are very large  $\geq 1$ Tb of information.
- Personalized medicine via genome sequencing is **HERE**.