



# Clinical and Research Genomics

## Spring 2022

### **Professor:**

Christopher E. Mason, Ph.D.

### **Instructors:**

Ebrahim Afshinnekoo, M.D.

Jaden Hastings, Ph.D.

### **TA:**

Chandrima Bhattacharya, M.S.

# Course Sessions:

- I. Sequencing Methods, Single-Cell Dynamics, and Molecular Detection Techniques (March 29<sup>th</sup>)
- II. RNA Sequencing, Epitranscriptomes, and Single Cell / Spatial Omics (April 5<sup>th</sup>)
- III. Epigenomes, DNA Modifications, and Chromatin Dynamics (April 12<sup>th</sup>)
- IV. Metagenomes, BGCs, and Metabolomics (April 19<sup>th</sup>)
- V. Complex Genome Re-arrangements, Transposons, and Tools for Genetic Variant Calling (April 26<sup>th</sup>)
- VI. Multi-Omics, Spatial Omics, and Machine Learning (May 3<sup>rd</sup>)
- VII. Synthetic Biology, Engineering Systems & Genome Ethics (May 10<sup>th</sup>)
- VIII. COVID-19 Tracking and Pathophysiology (May 17<sup>th</sup>)
- IX. Global Health and Beyond-Globe Health (Aerospace Medicine) (May 24<sup>th</sup>)

All classes also on Zoom:

<https://weillcornell.zoom.us/j/99565175034>

Meeting ID: 995 6517 5034 (Passcode: ClinGen22)

Course webpage:

<http://physiology.med.cornell.edu/faculty/mason/lab/clinicalgenomics/schedule.html>





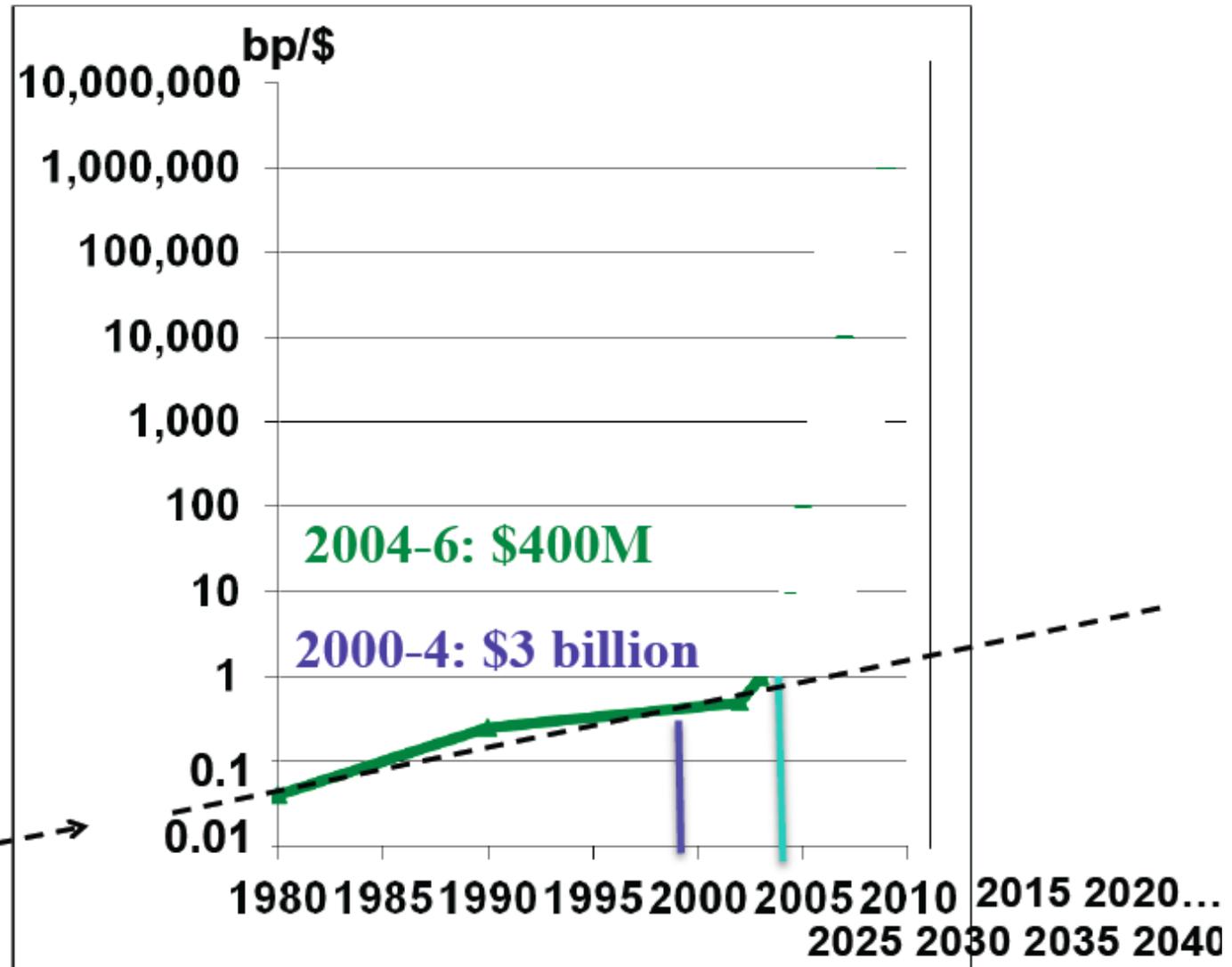
# Initially we expected a \$1K Genome in 2040

\$1000  
Genome

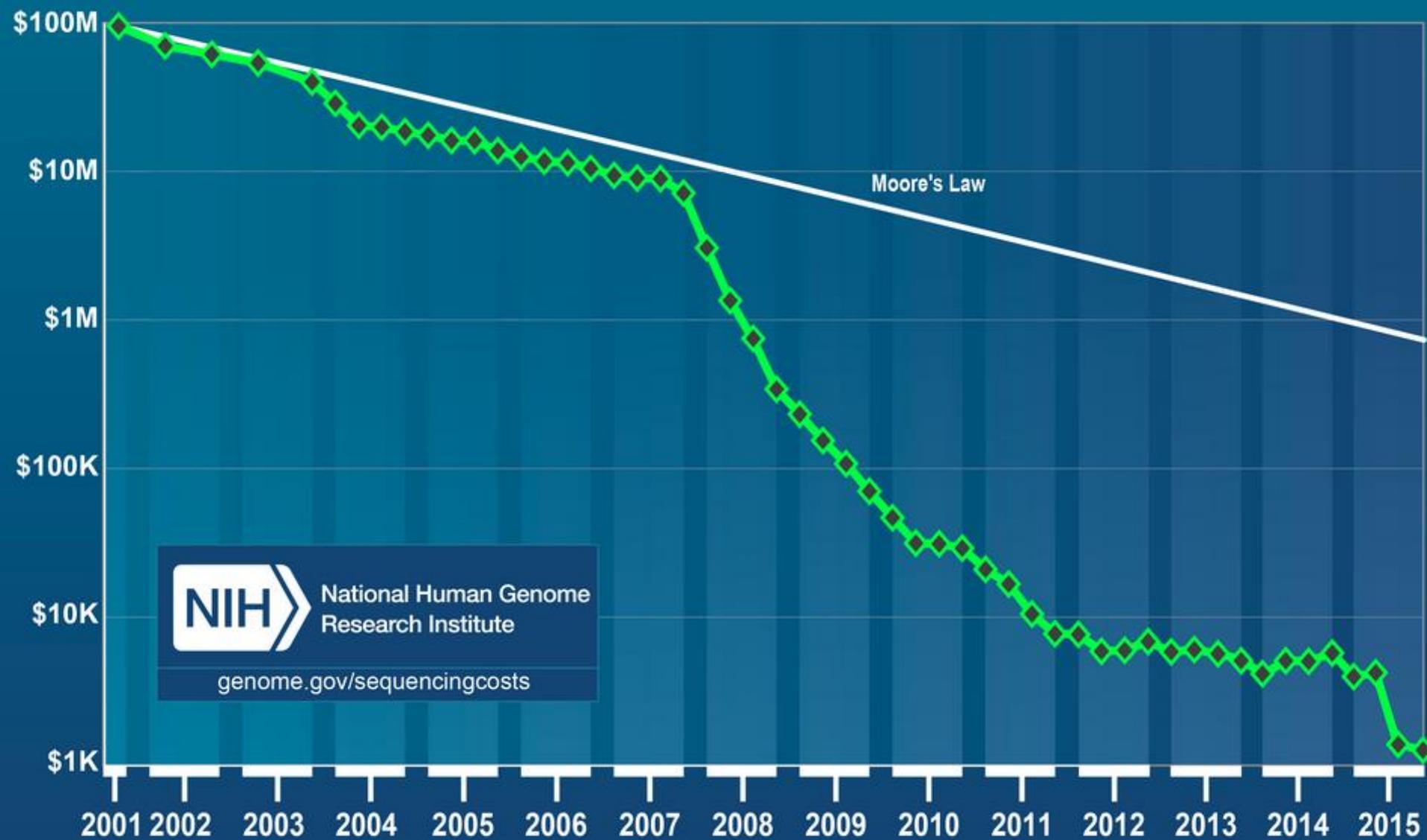
When?

2040

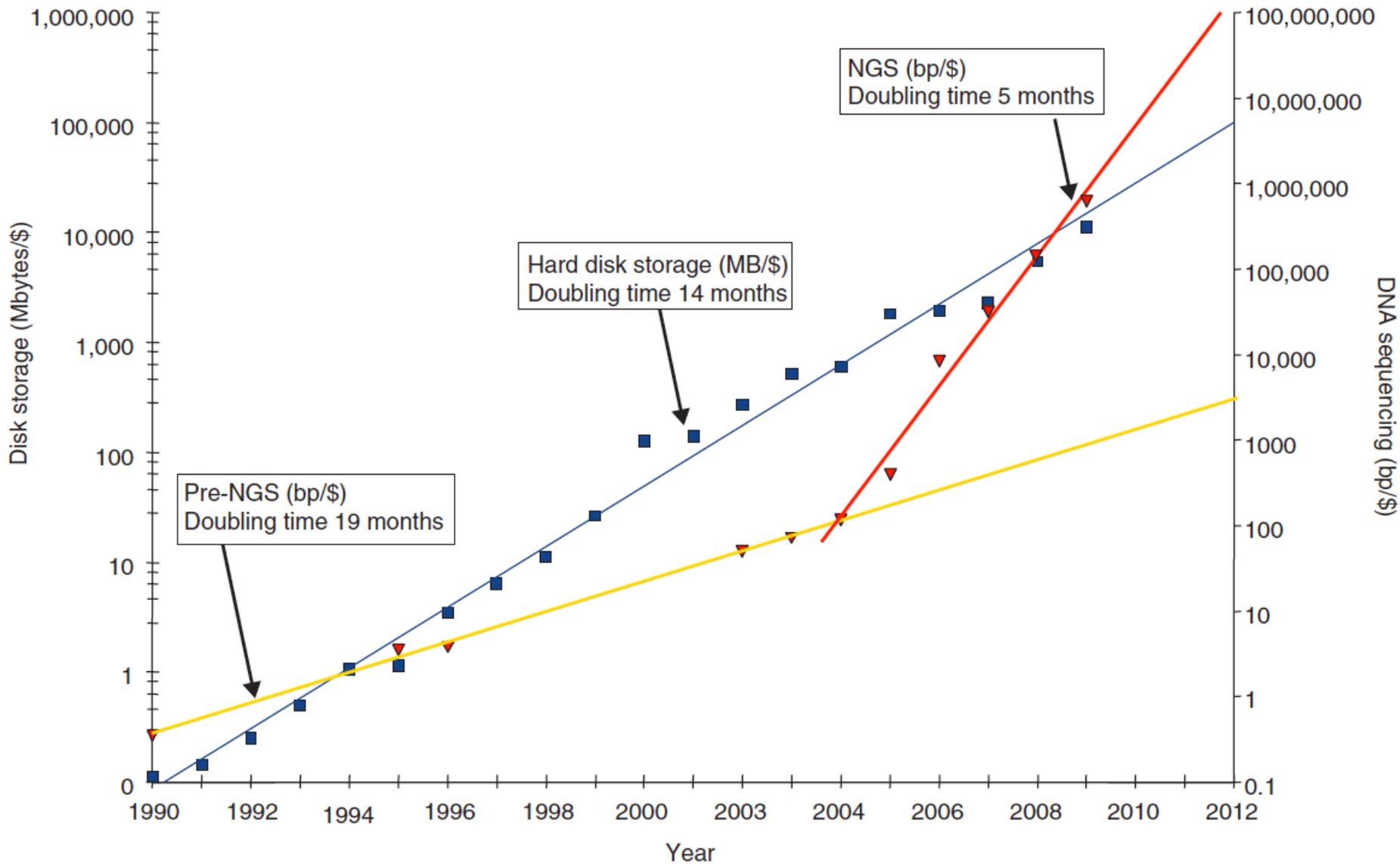
-----  
Moore's law  
1.5x/yr for  
electronics



# Cost per Genome



# More and more data





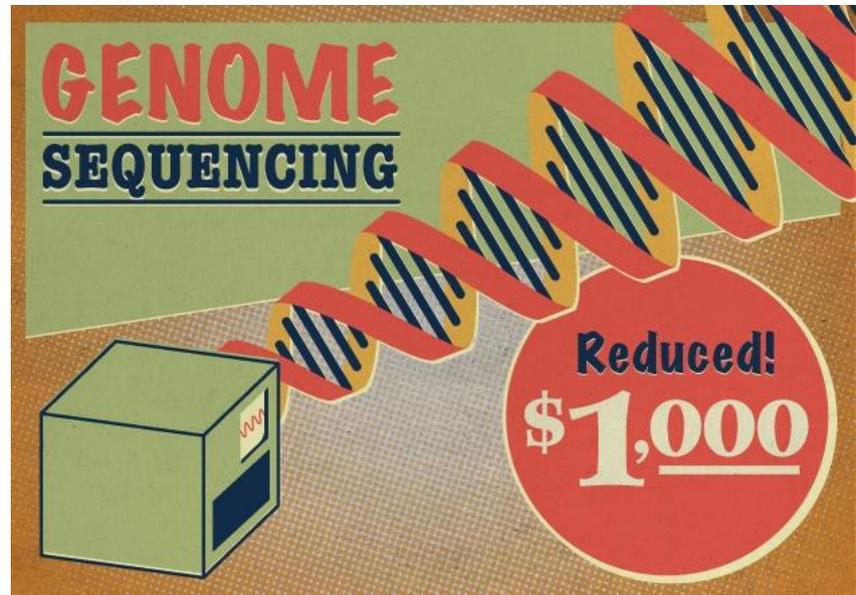


## Technology: The \$1,000 genome

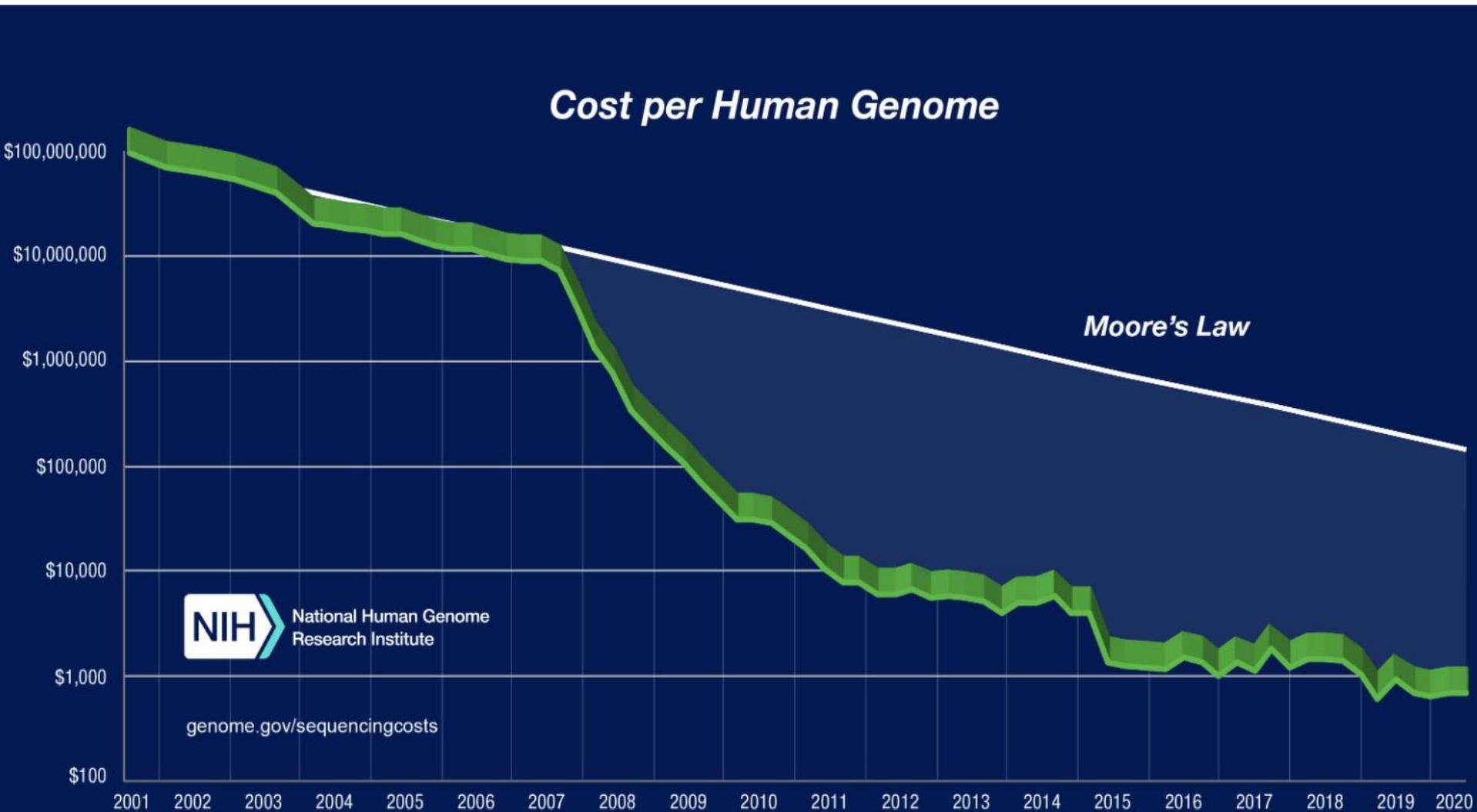
With a unique programme, the US government has managed to drive the cost of genome sequencing down towards a much-anticipated target.

Erika Check Hayden

19 March 2014



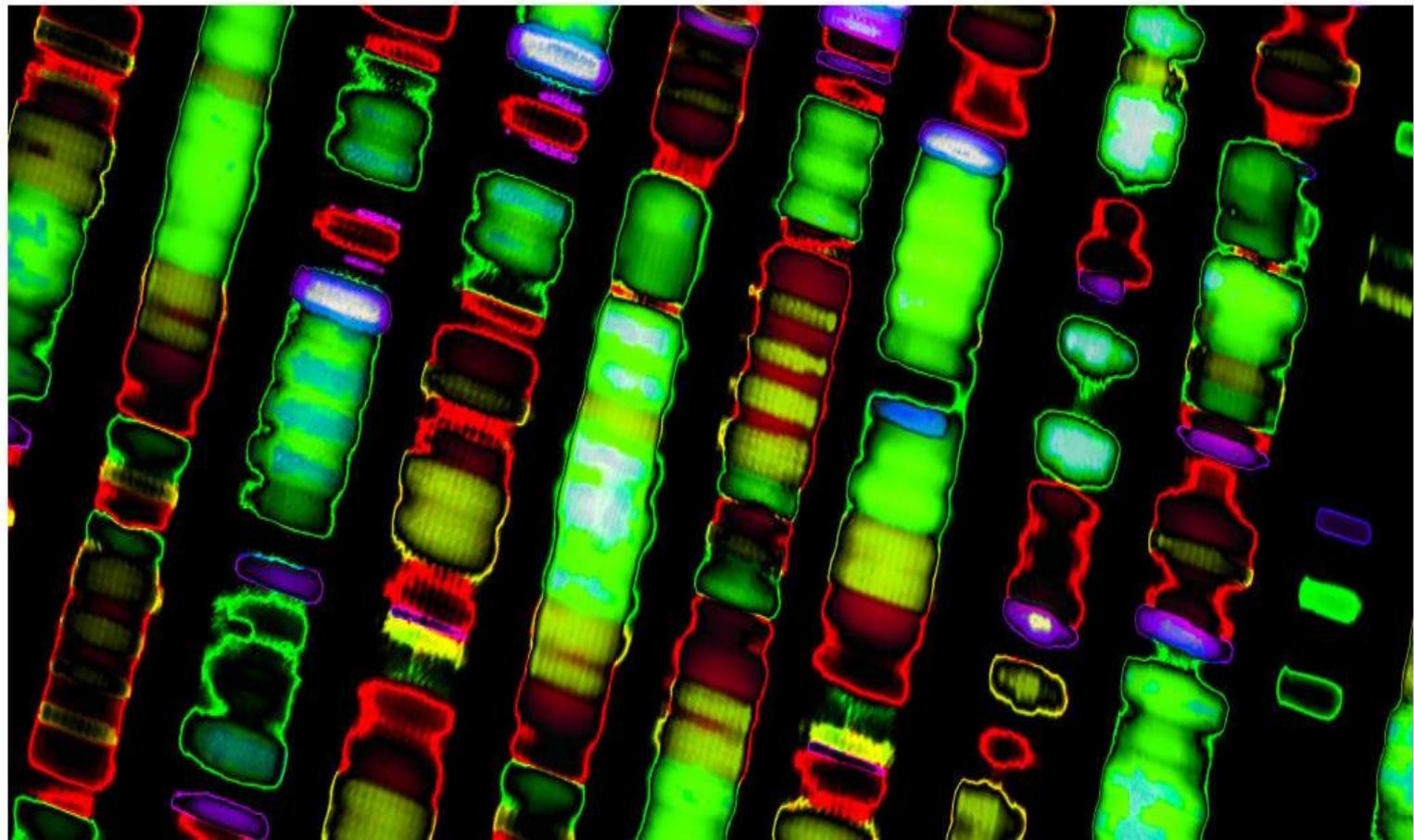
# Flatlined a little



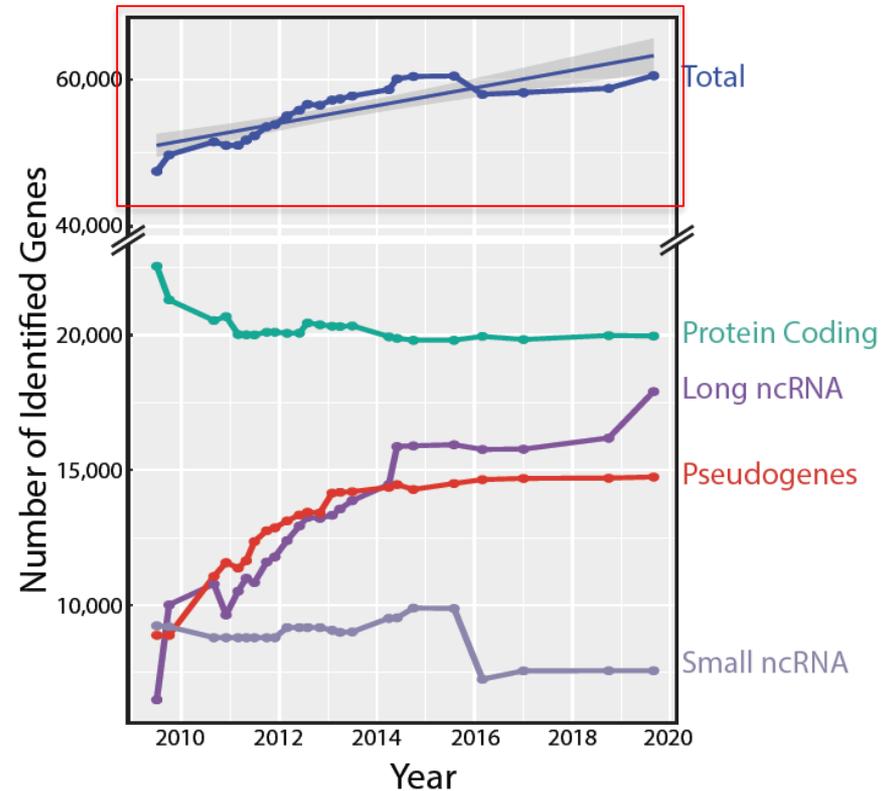
BUSINESS

# Illumina says it can deliver a \$100 genome — soon

- Twitter
- Facebook
- LinkedIn
- Email
- Share
- Print



# New human genes are still being found





### HUMAN

GENCODE 39 (09.12.21)



### MOUSE

GENCODE M28 (09.12.21)



### Tweets by @GencodeGenes

**GencodeGenes**  
@GencodeGenes

New single-nuclei RNA sequencing method finds interesting alt splicing expression in brain. Still plenty of novel transcripts out there for us to add to the GENCODE geneset! [nature.com/articles/s4158](https://www.nature.com/articles/s4158)



Mar 9, 2022

GencodeGenes Retweeted

**UCSC Genome Browser**  
@GenomeBrowser

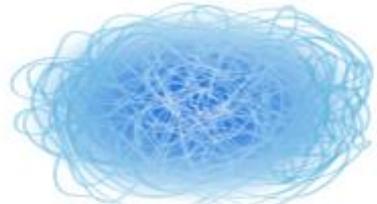
We have updated GENCODE Gene annotation tracks for human (V39 - hg19/hg38) and mouse (M28 - mm39) corresponding to Ensembl release 105.

The goal of the GENCODE project is to identify and classify all gene features in the human and mouse genomes with high accuracy based on biological evidence, and to release these annotations for the benefit of biomedical research and genome interpretation.

Every Day  
is the  
Best Day

# Human Genome Sequencing

Generating a Reference Genome Sequence  
(e.g., Human Genome Project)



Genomic DNA

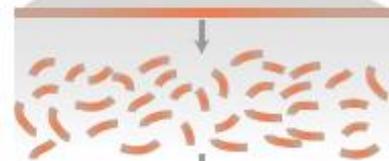
Break genome into large fragments and insert into clones



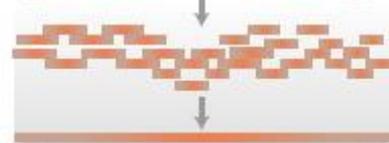
Order clones



Break individual clones into small pieces



Generate thousands of sequence reads and assemble sequence of clone

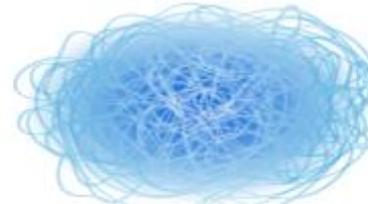


Assemble sequences of overlapping clones to establish reference sequence



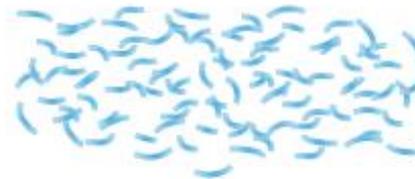
Reference Sequence

Generating a Person's Genome Sequence  
(e.g., Circa ~2016)



Genomic DNA

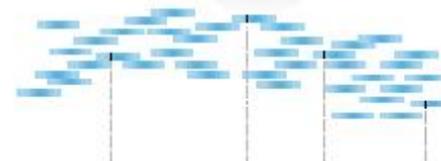
Break genome into small pieces



... TATGCGATGCGTATTTTCGTAA ...

Generate millions of sequence reads

Align sequence reads to established reference sequence

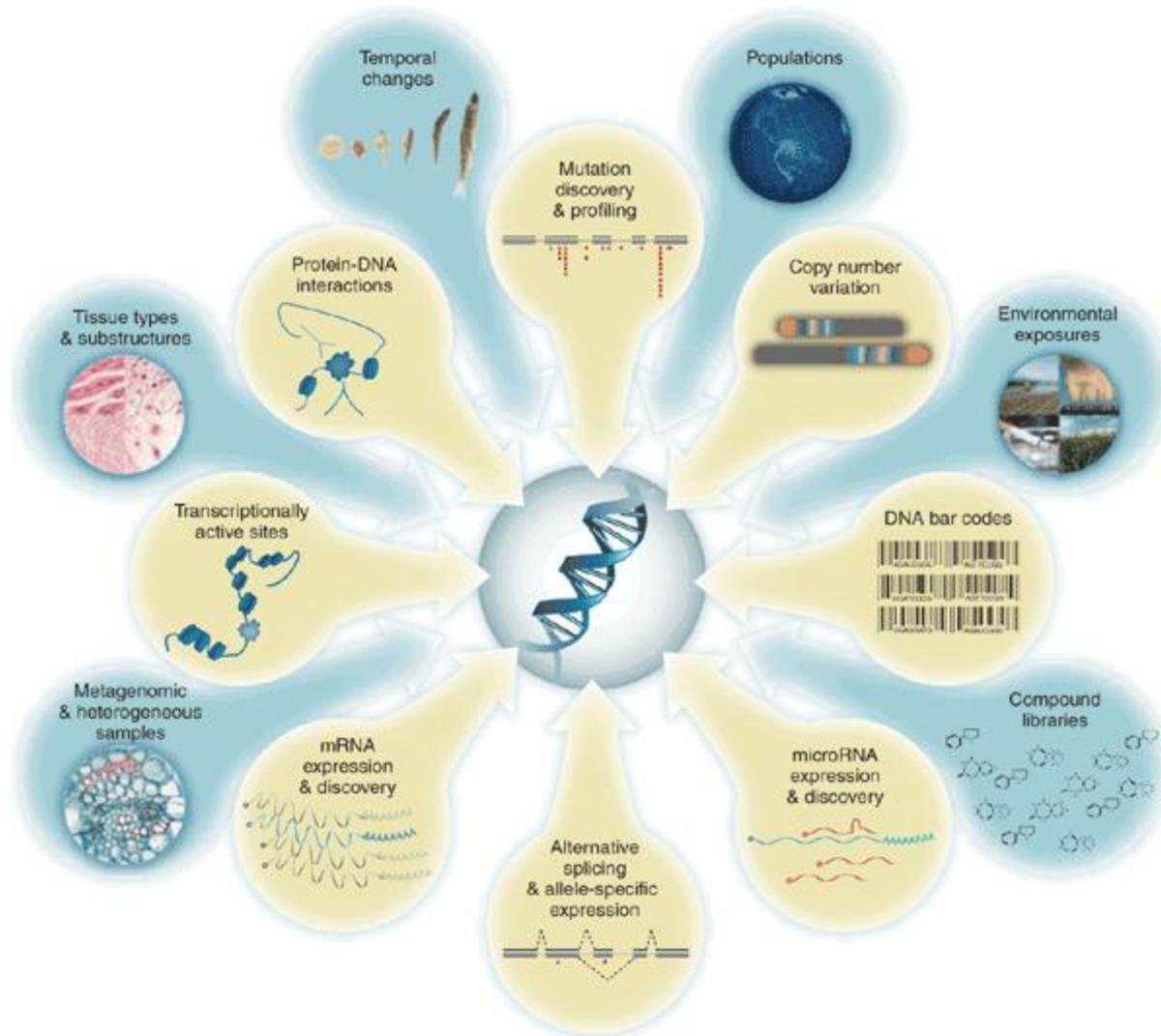


Reference Sequence

Deduce starting sequence and identify differences from reference sequence



Since DNA defines the biochemical recipe for the genesis of organisms, sequencing allows us to create molecular portraits of development and disease at single-base resolution.





[Declarations](#)

[References](#)

Musings | [Open Access](#)

## The \$1,000 genome, the \$100,000 analysis?

[Elaine R Mardis](#) 

*Genome Medicine* 2010 2:84

<https://doi.org/10.1186/gm205> | © BioMed Central Ltd 2010

**Published:** 26 November 2010



## Genomics England is delivering the **100,000 Genomes Project**.

We are creating a new genomic medicine service with the NHS – to support **better diagnosis and better treatments** for patients. We are also enabling medical research.

[More information about the 100,000 Genomes Project](#)

News story

### **Genome sequencing project reaches the halfway mark**

50,000 human genomes have now been sequenced from patients with cancer or rare diseases, under the 100,000 Genomes Project.

---

Published 28 February 2018



# ALL OF US<sup>SM</sup> RESEARCH PROGRAM

## All of Us Research Program

October 12, 2016

# PMI Cohort Program announces new name: the All of Us Research Program

The Precision Medicine Initiative® (PMI) Cohort Program will now be called the *All of Us* Research Program and will be the largest health and medical research program on precision medicine. A set of core values is guiding its development and implementation:

- Participation is open to all.
- Participants reflect the rich diversity of the U.S.
- Participants are partners.

[Scale and Scope](#)

[Participation](#)

[Program Components](#)

[Funding](#)

[FAQ](#)

[Advisory Groups](#)

[Events](#)

[Announcements](#)

[In the News](#)

[Multimedia](#)



# 1 million U.S. Veterans WGS



U.S. Department  
of Veterans Affairs



Search ORD

[VA SITE MAP \[A-Z\]](#)

[Health](#)

[Benefits](#)

[Burials & Memorials](#)

[About VA](#)

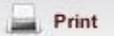
[Resources](#)

[News Room](#)

[Locations](#)

[Contact Us](#)

[VA](#) » [Health Care](#) » [Office of R&D](#) » [Mvp](#) » [Million Veteran Program \(MVP\)](#)



## Office of Research & Development

[ORD Home](#)

[About Us](#)

[Services](#)

[Programs](#)

[See All Programs](#)

[Animal Research](#)

[Biosafety & Biosecurity](#)

[Cooperative Studies Program \(CSP\)](#)

[Health Disparities & Minority Health](#)

[Million Veteran Program \(MVP\)](#)

### Million Veteran Program (MVP)

MVP is a national, **voluntary** research program funded entirely by the Department of Veterans Affairs Office of Research & Development. The goal of MVP is to partner with Veterans receiving their care in the VA Healthcare System to study how genes affect health. To do this, MVP will build one of the world's largest medical databases by safely collecting blood samples and health information from one million Veteran volunteers. Data collected from MVP will be stored anonymously for research on diseases like diabetes and cancer, and military-related illnesses, such as post-traumatic stress disorder. [Learn more.](#)



#### [Frequently Asked Questions](#)

- [How do I participate?](#)
- [Do I need to schedule an appointment to participate?](#)

Text size: [+](#) [-](#)

#### CONTACT MVP

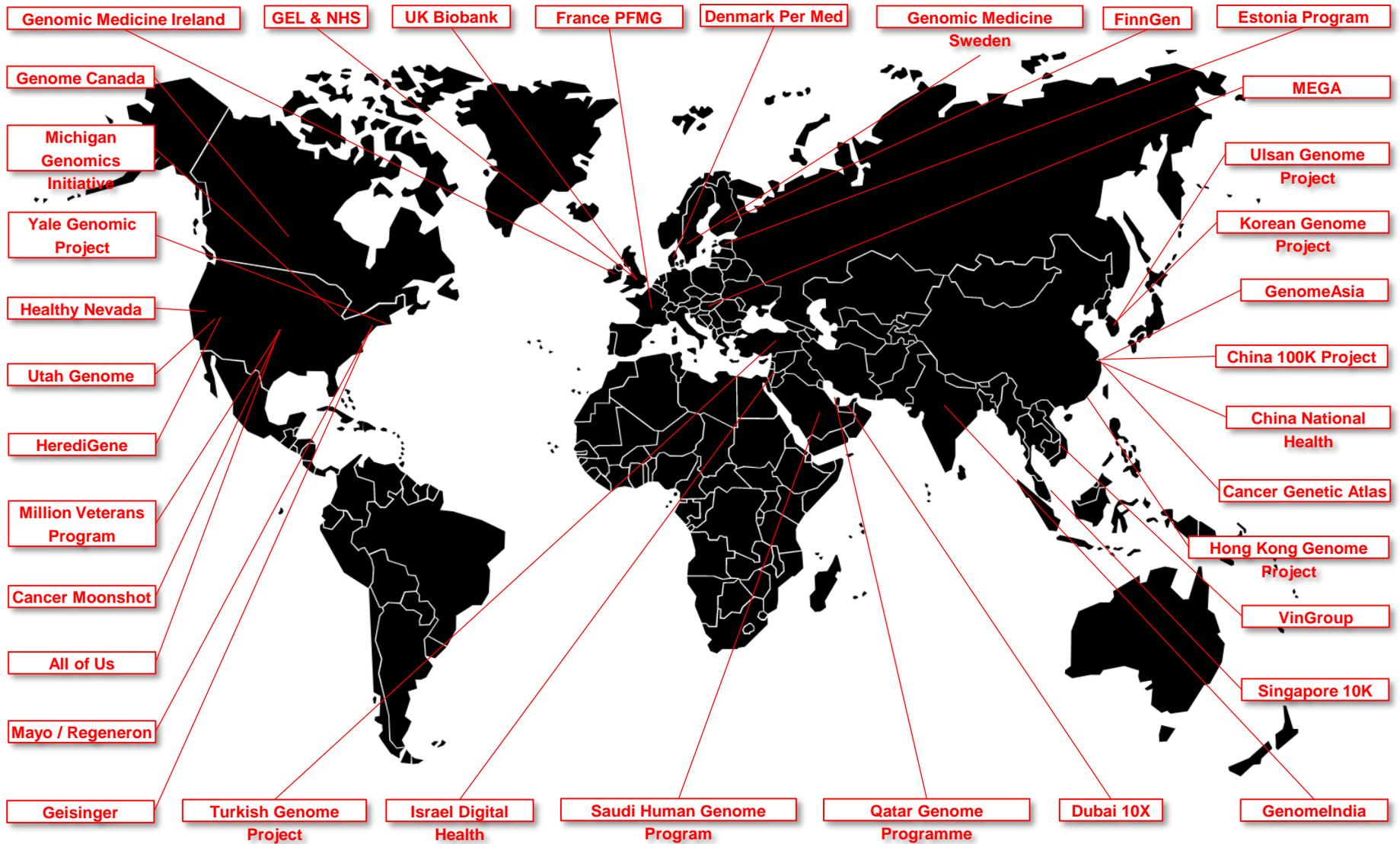
Contact the MVP Information Center toll-free at:

866-441-6075

#### INFORMED CONSENT



# POPULATION-SCALE NGS IS GLOBAL



# NHS to trial blood test to detect more than 50 forms of cancer

**Researchers hopes Galleri trial will be a 'gamechanger' for early diagnosis and save many lives**



▲ The Galleri blood test will be offered to 165,000 people in England from mid-2021, the vast majority of whom have no signs of the disease. Photograph: Jacqueline Larma/AP

Offered to 165,000 people in England from mid-2021 onward; no signs of disease.

Followed through 2023; If successful, move on to test 1M people in 2024-2025.

<https://www.theguardian.com/science/2020/nov/27/nhs-to-trial-blood-test-to-detect-more-than-50-forms-of-cancer>

# Specific genes can have significant impact

Myostatin (MSTN) homozygous nulls (-/-) give lean and large muscles

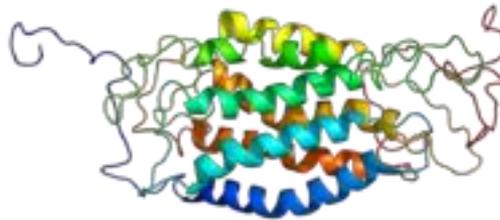


<http://thevoiceofnetizen.blogspot.com>

Low density lipoprotein receptor 5 (LRP5) heterozygotes (+/-) can have strong bones



C-C chemokine receptor type 5 (CCR5) homozygous nulls (-/-) have HIV protection



## Constitutional CHEK2 mutations are associated with a decreased risk of lung and laryngeal cancers

Cezary Cybulski\*, Bartłomiej Masojć, Dorota Oszutowska, Ewa Jaworowska<sup>1</sup>, Tomasz Grodzki<sup>2</sup>, Piotr Waloszczyk<sup>2</sup>, Piotr Serwatowski<sup>2</sup>, Juliusz Pankowski<sup>2</sup>, Tomasz Huzarski, Tomasz Byrski, Bohdan Górski, Anna Jakubowska, Tadeusz Dębniak, Dominika Wokołorczyk, Jacek Gronwald, Czesława Tarnowska<sup>1</sup>, Pablo Serrano-Fernández, Jan Lubiński and Steven A.Narod<sup>3</sup>

International Hereditary Cancer Center, Department of Genetics and Pathology, Pomeranian Medical University, ul. Połabska 4, 70-115 Szczecin, Poland, <sup>1</sup>Department of Otolaryngology and Laryngological Oncology, Pomeranian Medical University, ul. Unii Lubelskiej, 71–252 Szczecin, Poland, <sup>2</sup>Lung Diseases Hospital, ul. Sokołowskiego 11, 70–891 Szczecin, Poland and <sup>3</sup>Women's College Research Institute, Toronto, Ontario M5G 1N8, Canada

\*To whom correspondence should be addressed. Tel: +48 91 466 1532;  
Fax: +48 91 466 1533;  
Email: cezarycy@sci.pam.szczecin.pl

**Mutations in the CHEK2 gene have been associated with increased risks of breast, prostate and colon cancer. In contrast, a previous report suggests that individuals with the I157T missense variant of the CHEK2 gene might be at decreased risk of lung cancer and upper aero-digestive cancers. To confirm this hypothesis, we genotyped 895 cases of lung cancer, 430 cases of laryngeal cancer and 6391 controls from Poland for four founder alleles in the CHEK2 gene, each of which has been associated with an increased risk of cancer at several sites. The presence of a CHEK2 mutation was protective against both lung cancer [odds ratio (OR) = 0.3; 95% confidence interval (CI) 0.2–0.5;  $P = 3 \times 10^{-8}$ ] and laryngeal cancer (OR = 0.6; 95% CI 0.3–0.99;  $P = 0.05$ ). The basis of the protective effect is unknown, but may relate to the reduced viability of lung cancer cells with a CHEK2 mutation. Lung cancers frequently possess other defects in genes in the DNA damage response pathway (e.g. p53 mutations) and have a high level of genotoxic DNA damage induced by tobacco smoke. We speculate that lung cancer cells with impaired CHEK2 function undergo increased rates of cell death.**

### Introduction

Germ line mutations in CHEK2 have been associated with a range of cancer types, in particular of the breast and the prostate, but cancers of

of Brennan *et al.* We have extended our series of lung cancer cases from 272 to 895 and our control sample from 4000 to 6391. We have also identified a fourth deleterious CHEK2 allele (a large deletion of exons 9 and 10). Because smoking is the principal risk factor for lung cancer in Poland and elsewhere, we asked whether the protective effect of CHEK2 might extend to laryngeal cancer patients as well.

### Materials and methods

We studied 895 unselected cases of lung cancer (226 women and 669 men) diagnosed in the Lung Diseases Hospital in Szczecin, Poland, between 2004 and 2006. We also ascertained 430 consecutive, unselected patients with squamous cell carcinoma of the larynx (70 women and 360 men) at Department of Otolaryngology and Laryngological Oncology of the Pomeranian Medical University, Szczecin, Poland, during the period 2001–2004. Patients were recruited from the oncology services of the contributing hospitals and were unselected for age or family history. Patients were approached by a member of the study team during an outpatient visit to the oncology clinic and were asked if they wished to participate. Patient acceptance rates exceeded 80% for both cancer sites. Patients provided written informed consent. A blood sample of 10 cc was then drawn for DNA extraction. Two hundred and seventy-two of the lung cancer patients have been included in our previous study (5). The mean age of diagnosis of the lung cancer patients was 61.4 years (range 29–88 years) and of the laryngeal cancer patients was 58.2 years (range 30–84). Patients completed a questionnaire about their smoking habits at the time of cancer diagnosis. Smoking histories were available for 818 of 895 (91%) lung cancer cases and for 387 of 430 (90%) laryngeal cancer cases. The study was approved by the Ethics Committee of the Pomeranian Medical University in Szczecin.

### Unmatched analysis

In the unmatched analysis, four non-overlapping control groups were combined in order to maximize the number of controls.

The first control group of 1896 healthy adults, including 1079 women (age range 15–91, mean 58.3) and 817 men (age range 23–90, mean 59.4). These controls were selected at random from the computerized patient lists of five large family practices located in the region of Szczecin. These healthy adults were invited to participate by mail and participated in 2003 and 2004. Participation rates for this group exceeded 70%. During the interview, the goals of the study were explained, informed consent was obtained, genetic counselling was given and a blood sample was taken for DNA analysis. A detailed family history of cancer was taken (first- and second-degree relatives included). Probands were included regardless of their cancer family history status. Individuals affected with any malignancy were excluded from the study.

The second control group consisted of 1417 unselected young adults (705 women and 712 men; age range 18–35, mean 24.3) from Szczecin metropolitan region who submitted a blood sample for paternity testing between 1994 and 2001.

The third control group consisted of 2183 children from nine cities in Poland



Article | [Open Access](#) | Published: 03 October 2019

# Towards precision medicine: interrogating the human genome to identify drug pathways associated with potentially functional, population- differentiated polymorphisms

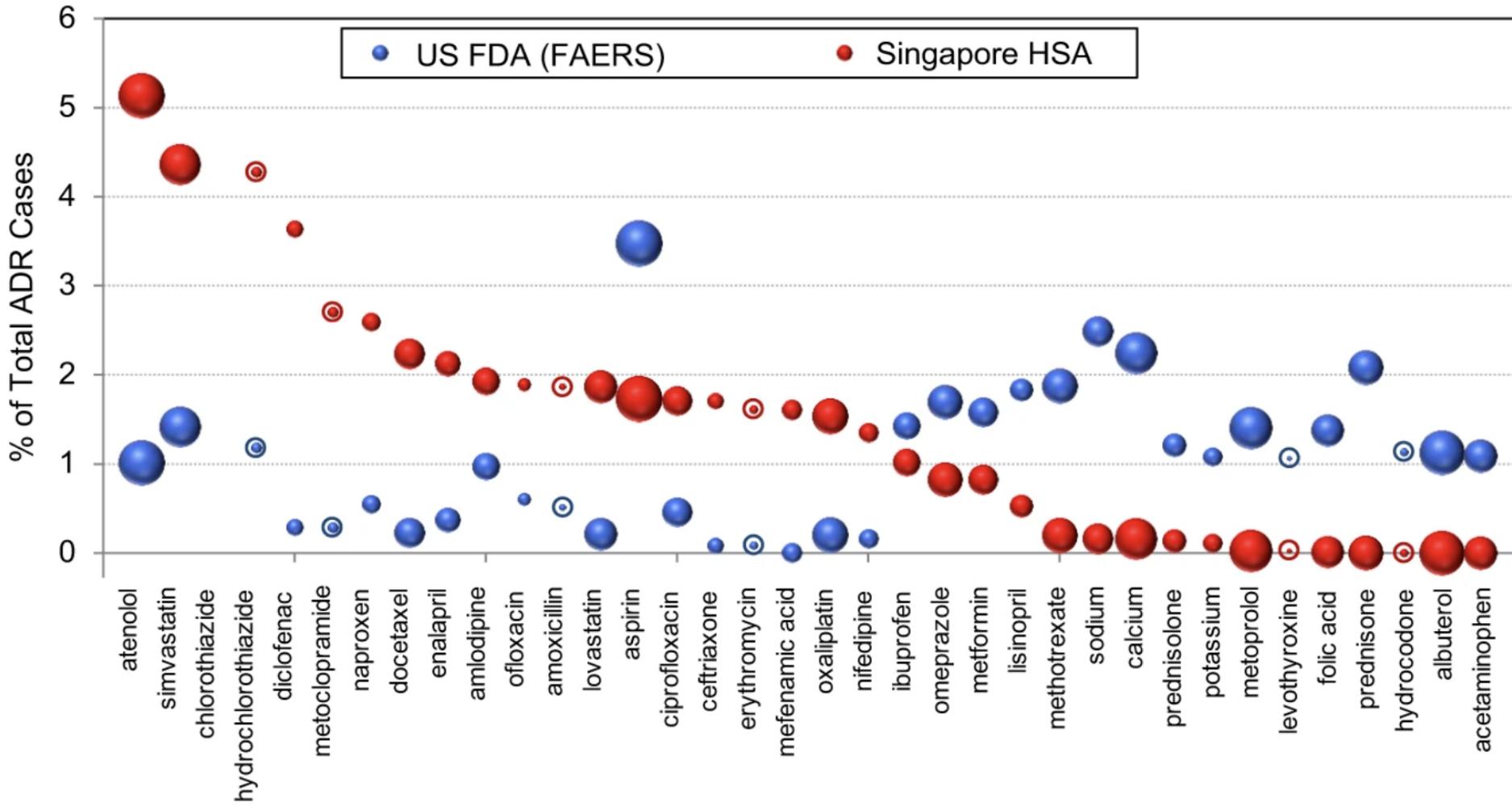
Maulana Bachtiar, Brandon Nick Sern Ooi, Jingbo Wang, Yu Jin, Tin Wee Tan, Samuel S. Chong & Caroline G. L. Lee 

*The Pharmacogenomics Journal* (2019) | [Download Citation](#) ↓

7 Accesses | **22** Altmetric | [Metrics](#) >>

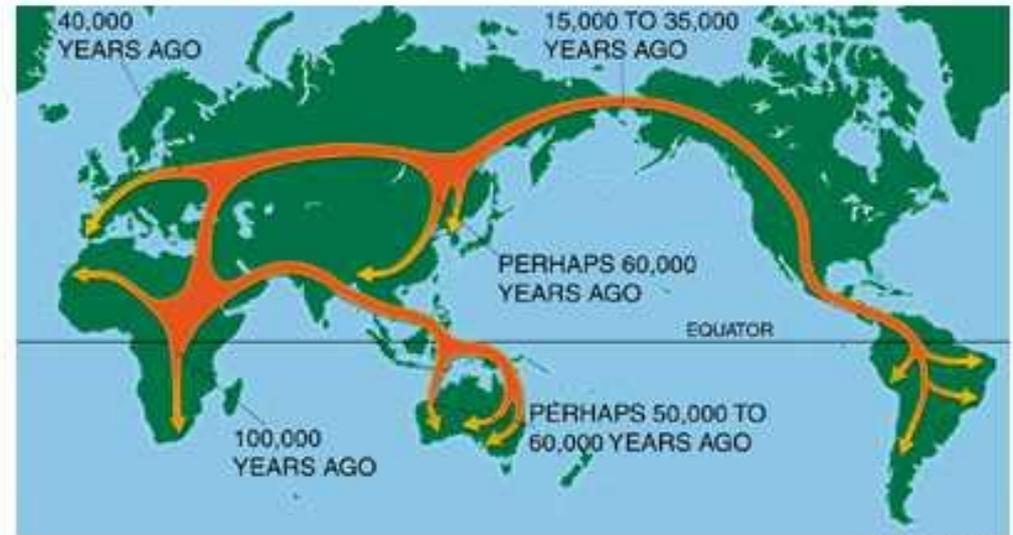
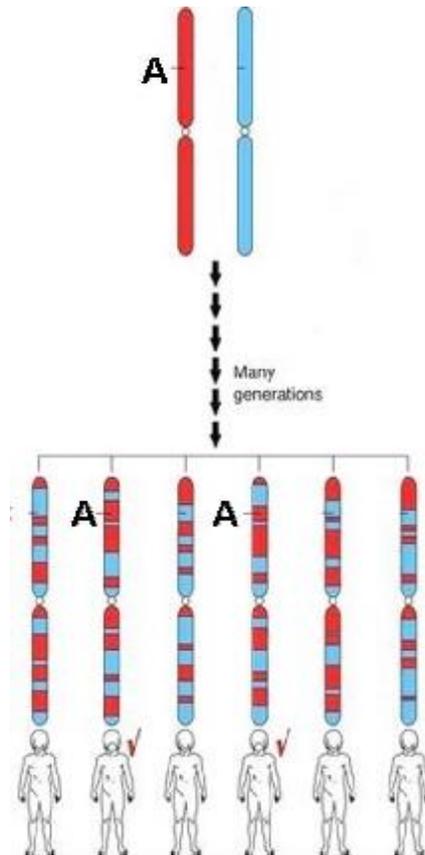
<https://www.nature.com/articles/s41397-019-0096-y>

## Top 20 suspected ADR drugs reported to Singapore HSA and US FDA



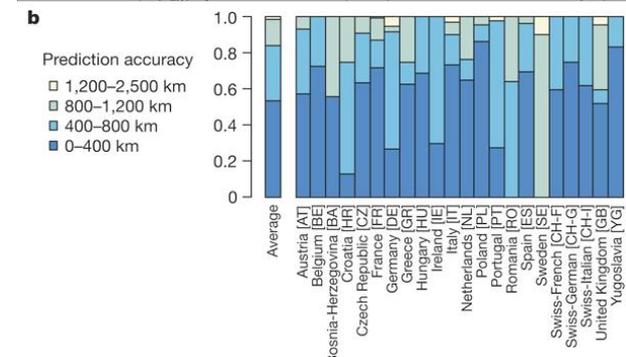
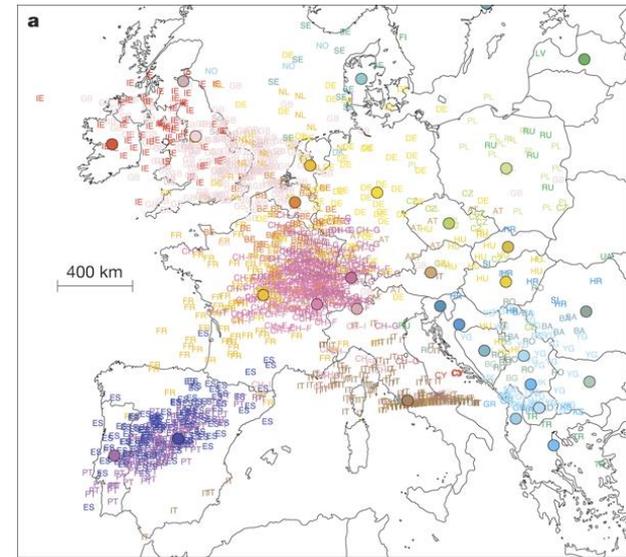
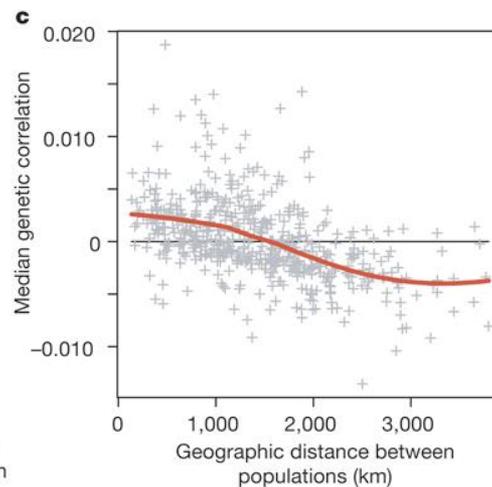
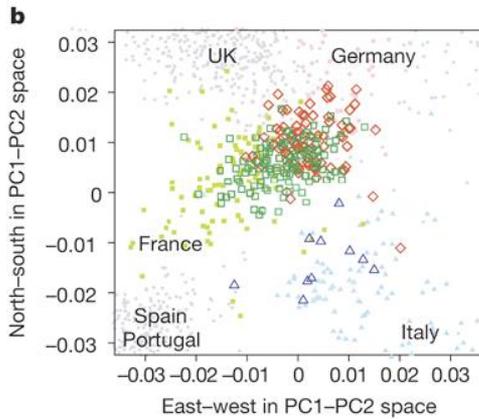
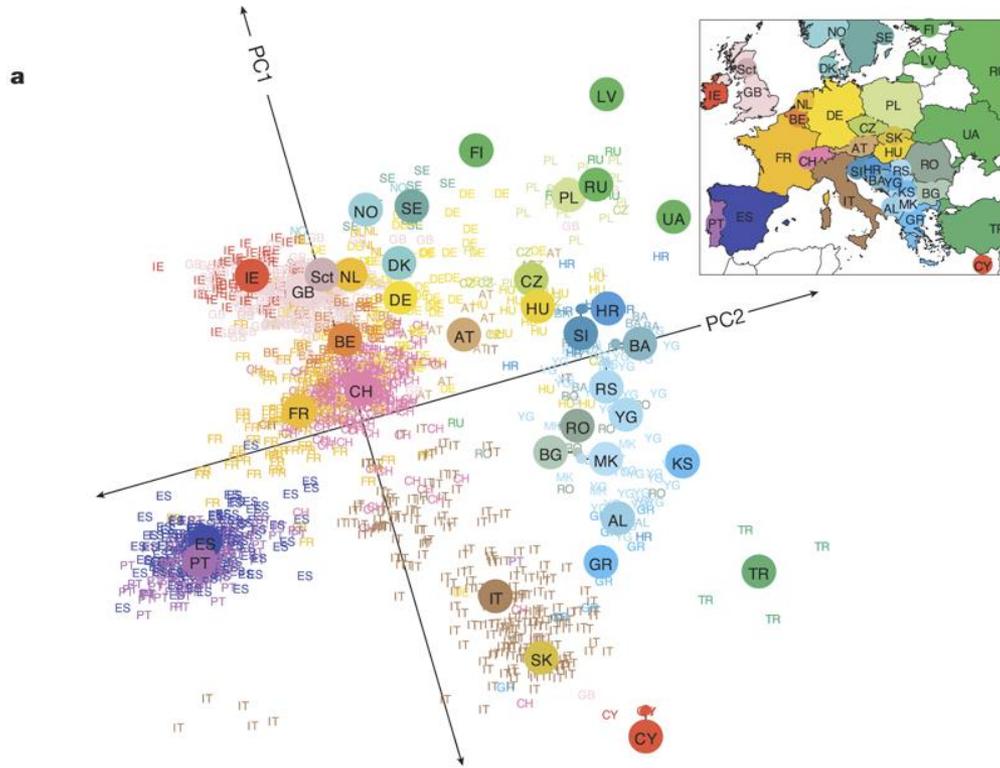
<https://www.nature.com/articles/s41397-019-0096-y>

# Our genes come from the migration patterns of haplotypes throughout human history (“Population Stratification”)



Tom Moore

# Genotype data can even predict your birthplace



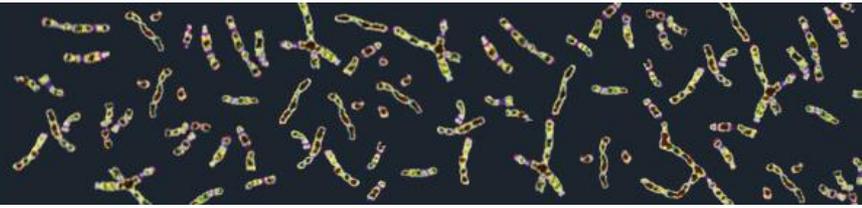
- French-speaking Swiss
- ◇ German-speaking Swiss
- △ Italian-speaking Swiss
- French
- German
- Italian

Genes mirror geography within Europe  
 Novembre *et al.*, 2008

# Large impact for normal genomes and diseases, especially cancer

## 1000 Genomes

A Deep Catalog of Human Genetic Variation



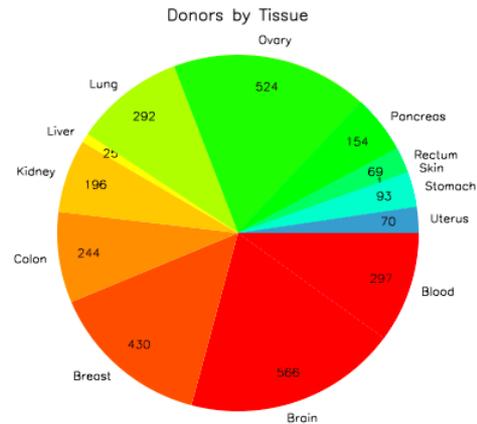
The Cancer Genome Atlas  
Data Portal



*Understanding genomics  
to improve cancer care*

ICGC DATASET VERSION 8 (MARCH 15TH, 2012)

Cancer Projects: 29



Total Donors: 3,561



International  
Cancer Genome  
Consortium

ICGC Goal: To obtain a comprehensive description of genomic, epigenomic, and transcriptomic (GET) changes in 50 different tumor types and/or subtypes which are of clinical and societal importance across the globe.

The cBio Cancer Genomics Portal provides **visualization, analysis** and **download** of large-scale **cancer genomics** data sets.

Please adhere to [the TCGA publication guidelines](#) when using any TCGA data in your

Filtered in 66 (48%) of cases.

Total 66 cases with alter  
altered

## Data Sets

The Portal contains data for **10410 tumor samples from 31 cancer studies.** [Details.]

Home
Query the Data
Download Data
Tools
About the Data
Publication Guidelines

### Home

---

## TCGA Data Portal Overview

We provide 3 ways to download data: The Cancer Genome Atlas (TCGA) Data Portal provides a platform for researchers to search, download, and analyze data sets generated by TCGA. It contains clinical information, genomic characterization data, and high-throughput sequencing analysis of the tumor genomes.

The TCGA Data Portal does not host lower levels of sequence data. NCI's [Cancer Genomics Hub \(CGHub\)](#) is the new secure repository for storing, cataloging, and accessing sequence related data. New users must still apply for authorized access through NCBI's [Database of Genotypes and Phenotypes \(dbGaP\)](#).

Query the Data >

Search summarized data for genes, patients and pathways

Download Data >

Choose from three ways to download data

Available Cancer Types	# Patients with Samples	# Downloadable Tumor Samples	Date Last Updated (mm/dd/yy)
<a href="#">Acute Myeloid Leukemia [LAML]</a>	202	200	02/15/13
<a href="#">Bladder Urothelial Carcinoma [BLCA]</a>	171	153	03/07/13
<a href="#">Brain Lower Grade Glioma [LGG]</a>	232	222	03/08/13
<a href="#">Breast invasive carcinoma [BRCA]</a>	956	940	03/08/13

### Announcements

---

**03/06/2013 - DCC Software Released**

The software release scheduled for today has been successfully completed and the TCGA Data Portal has been returned to operation. A complete list of the issues addressed in this release can be found on the TCGA Wiki [release notes](#) and for those with JIRA access the tickets covered in this release can be found on the wiki [here](#). Please note the release notes have been updated since they were published.

If you have any questions or concerns about this release, contact [tcga-dcc-binf-l@list.nih.gov](mailto:tcga-dcc-binf-l@list.nih.gov).

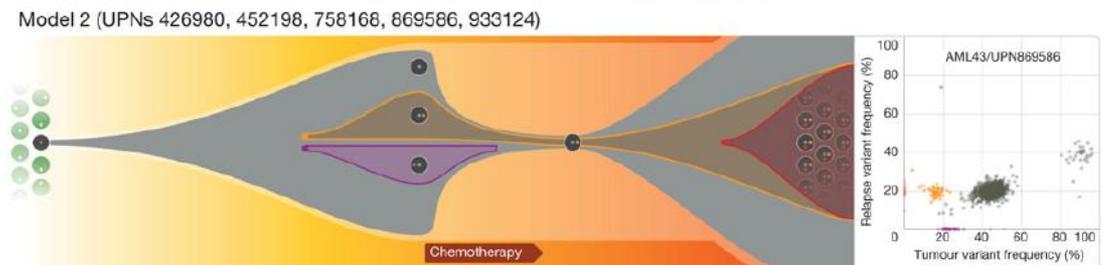
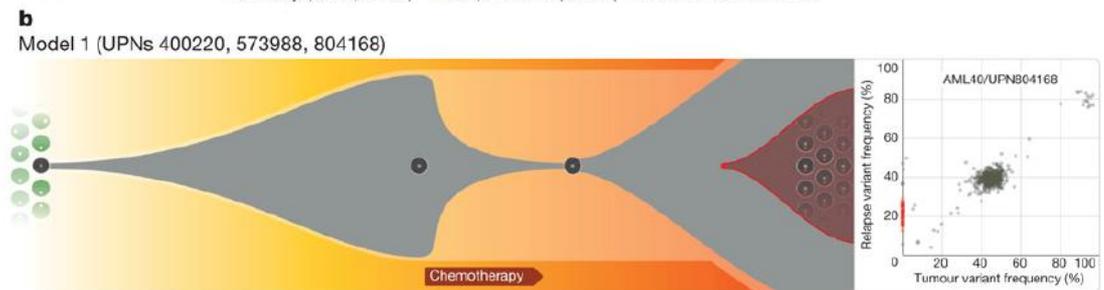
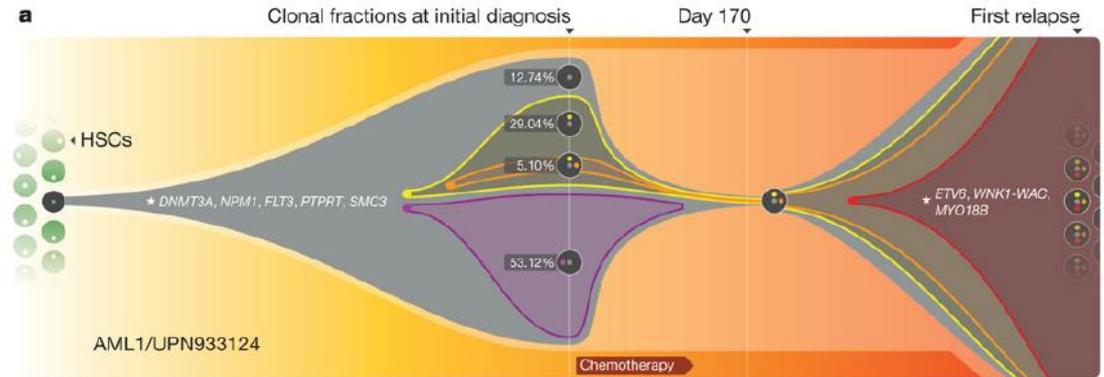
---

**02/25/2013 - DCC Software Released**

The software release scheduled for today has been successfully completed and the TCGA Data Portal has been returned to operation. A complete list of the issues addressed in this release can be found on the TCGA Wiki [Release Notes](#) and for those with JIRA access the tickets covered in this release can be found on the wiki [here](#)

If you have any questions or concerns about this release, contact [tcga-dcc-](mailto:tcga-dcc-)

# We can also observe the dynamics and evolution of cancers



**Cancer Genomics and Chemo**

Research from The Genome Institute and colleagues suggests chemotherapy may contribute to relapse in some patients with acute myeloid leukemia.

More on genomics and chemo >>

AML Case 803124 Tumor Evolution

ETV6, WNK1-WAC, MVO18B

DNMT3A, NPM1, FLT3, PTPRT, SMC3

Relapse

Ding L, et.al, Clonal evolution in relapsed acute myeloid leukemia revealed by whole-genome sequencing. Nature. 2012 Jan 11;481(7382):506-10.

# And look beyond just humans

## Genome 10K Project

To understand how complex animal life evolved through changes in DNA and use this knowledge to become better stewards of the planet



The Genome 10K project: Assembling a "Noah's Ark" of genomic data to save dying species.



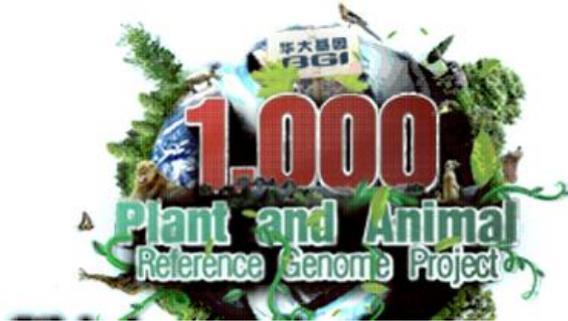
<https://genome10k.soe.ucsc.edu/>

<https://www.hgsc.bcm.edu/i5k-pilot-project-summary>

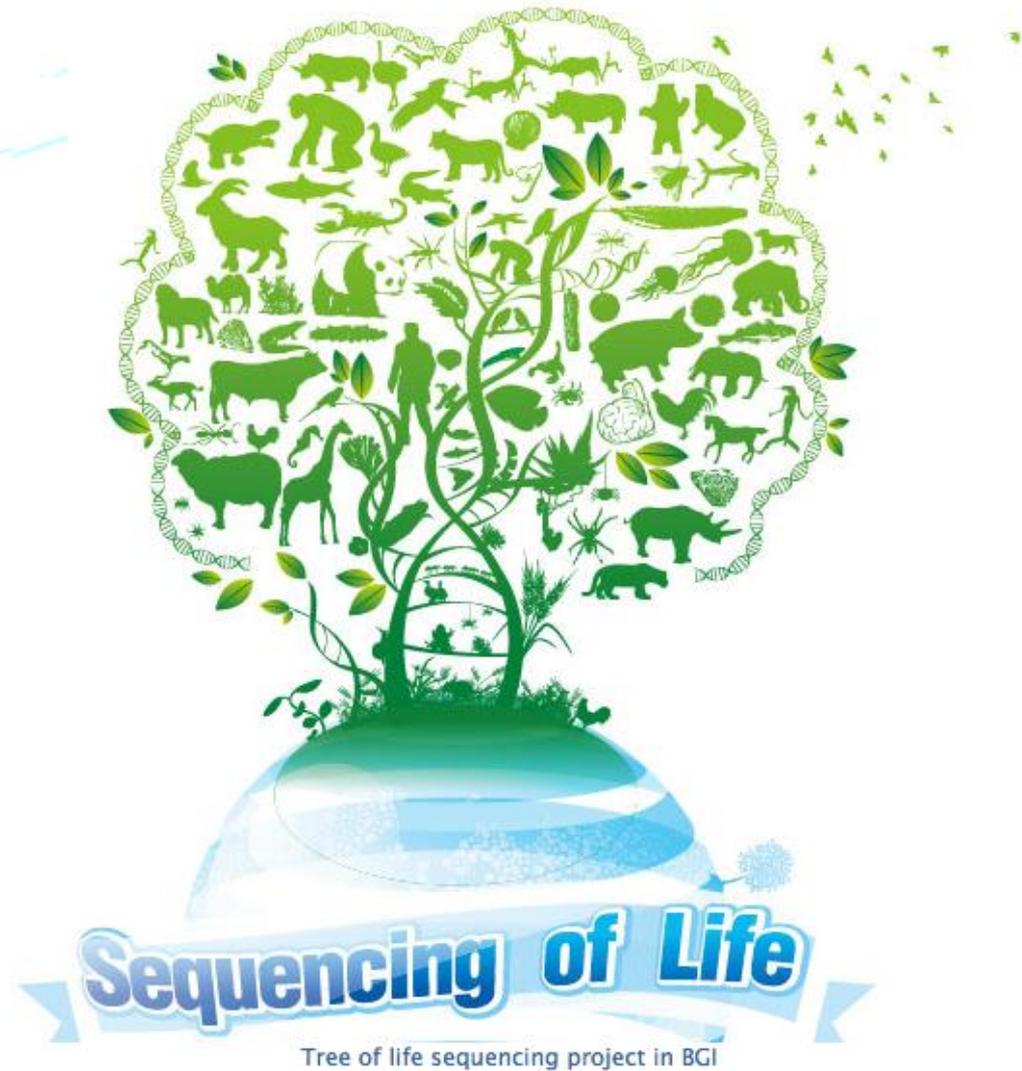




# Plants as well!



华大基因  
BGI



<http://idl.genomics.cn/page/pa-research.jsp>

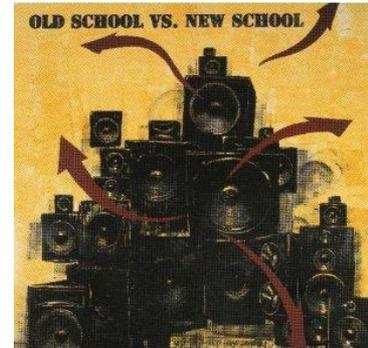
The Tech

# Sequencing Technologies

1. “Old School” dye-terminator sequencing (Sanger). 300-1000bp

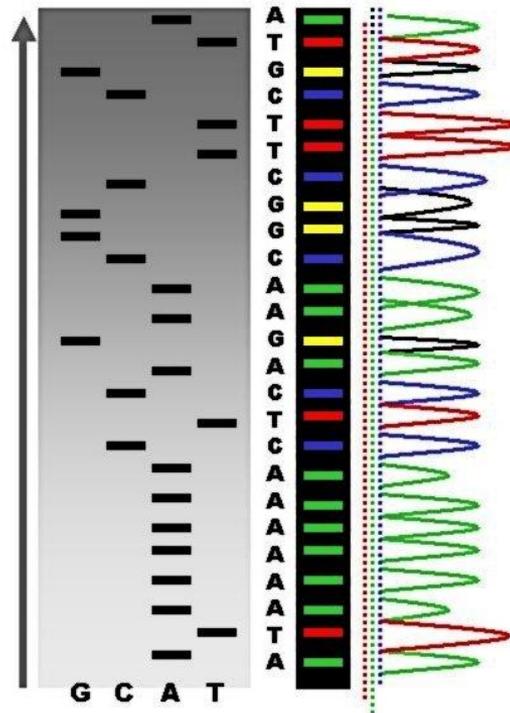
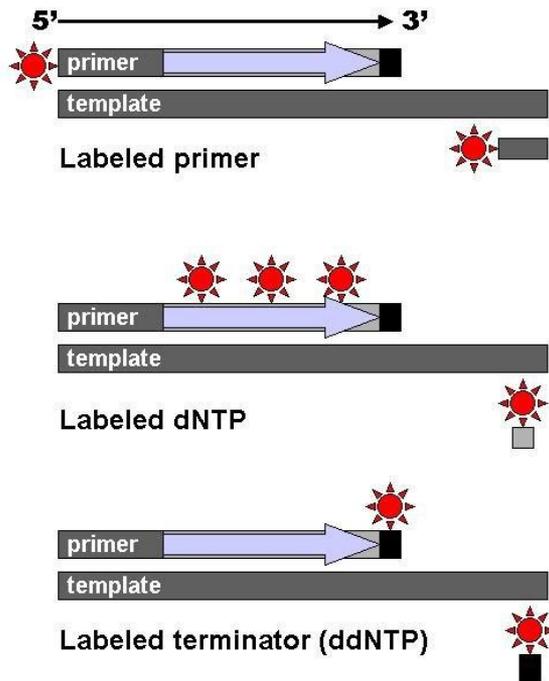
2. “New School” methods

- a. Emulsion PCR Pyrosequencing
- b. Solid-phase amplification sequencing by synthesis (clonal or single molecule)
- c. Sequencing by ligation
- d. Single-molecule, real-time (SMRT) sequencing
- e. Electrical sequencing



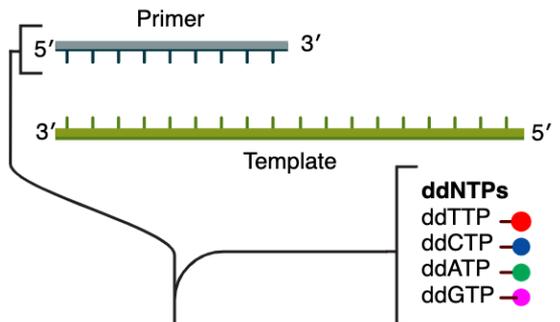
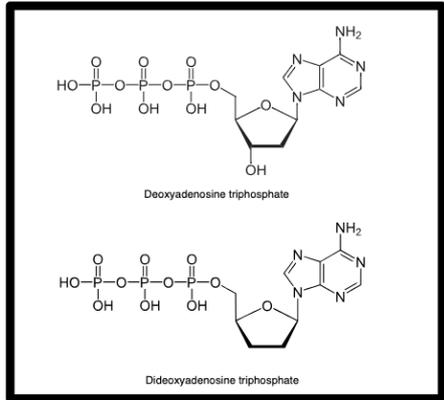
# Sequencing Technologies

## 1. “Old School” dye-terminator sequencing (Sanger). 300-1000bp



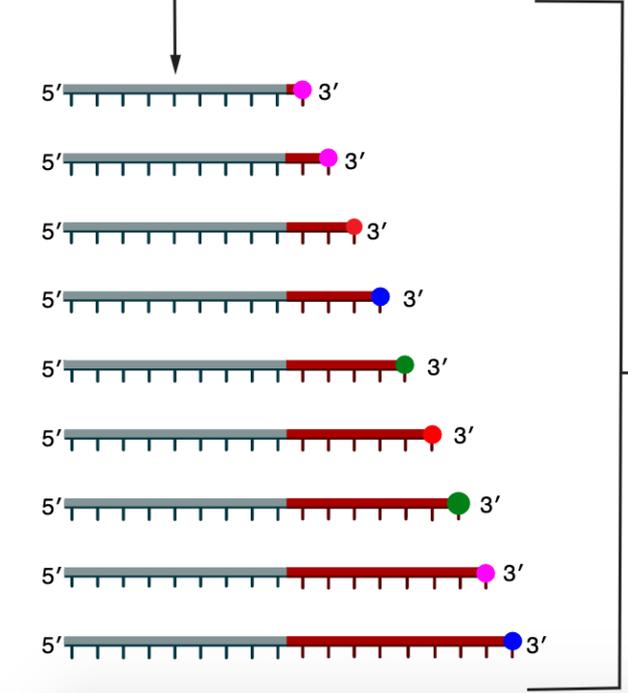
**① Reaction mixture**

- ▶ **Primer and DNA template** ▶ **DNA polymerase**
- ▶ **ddNTPs with flouochromes** ▶ **dNTPs (dATP, dCTP, dGTP, and dTTP)**

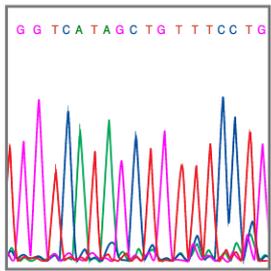
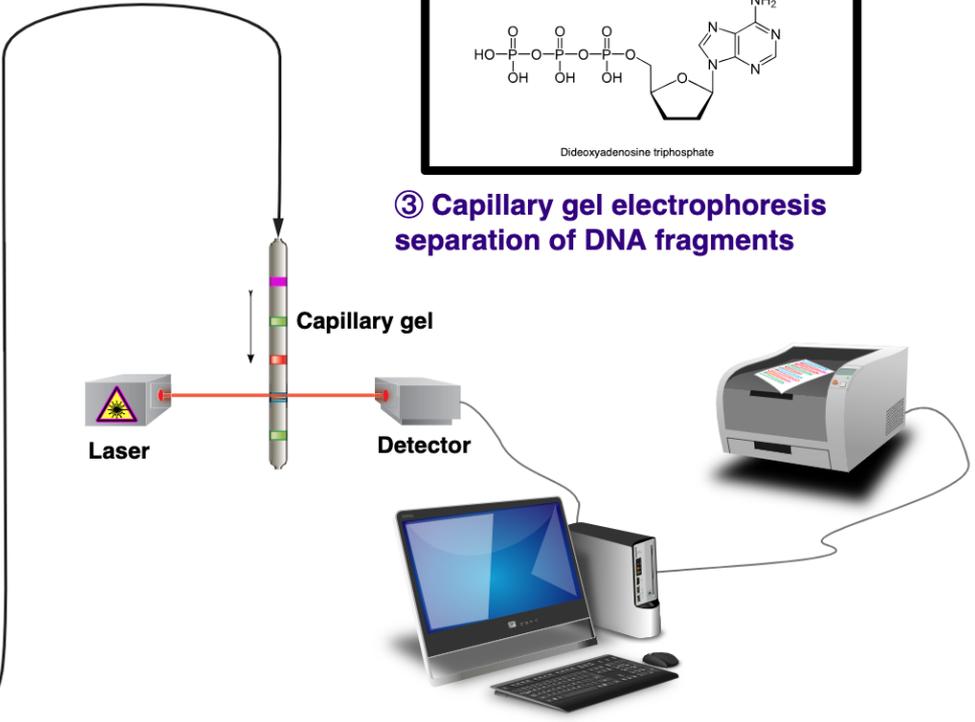


- ddNTPs**
- ddTTP ●
  - ddCTP ●
  - ddATP ●
  - ddGTP ●

**② Primer elongation and chain termination**



**③ Capillary gel electrophoresis separation of DNA fragments**



**Chromatograph**

**④ Laser detection of flouochromes and computational sequence analysis**

# By 2009, many MPS options emerged

Platform	Library/ template preparation	NGS chemistry	Read length (bases)	Run time (days)	Gb per run	Machine cost (US\$)	Pros	Cons	Biological applications	Refs
Roche/454's GS FLX Titanium	Frag, MP/ emPCR	PS	330*	0.35	0.45	500,000	Longer reads improve mapping in repetitive regions; fast run times	High reagent cost; high error rates in homo- polymer repeats	Bacterial and insect genome <i>de novo</i> assemblies; medium scale (<3 Mb) exome capture; 16S in metagenomics	D. Muzny, pers. comm.
Illumina/ Solexa's GA <sub>II</sub>	Frag, MP/ solid-phase	RTs	75 or 100	4 <sup>†</sup> , 9 <sup>§</sup>	18 <sup>†</sup> , 35 <sup>§</sup>	540,000	Currently the most widely used platform in the field	Low multiplexing capability of samples	Variant discovery by whole-genome resequencing or whole-exome capture; gene discovery in metagenomics	D. Muzny, pers. comm.
Life/APG's SOLiD 3	Frag, MP/ emPCR	Cleavable probe SBL	50	7 <sup>†</sup> , 14 <sup>§</sup>	30 <sup>†</sup> , 50 <sup>§</sup>	595,000	Two-base encoding provides inherent error correction	Long run times	Variant discovery by whole-genome resequencing or whole-exome capture; gene discovery in metagenomics	D. Muzny, pers. comm.
Polonator G.007	MP only/ emPCR	Non- cleavable probe SBL	26	5 <sup>§</sup>	12 <sup>§</sup>	170,000	Least expensive platform; open source to adapt alternative NGS chemistries	Users are required to maintain and quality control reagents; shortest NGS read lengths	Bacterial genome resequencing for variant discovery	J. Edwards, pers. comm.
Helicos BioSciences HeliScope	Frag, MP/ single molecule	RTs	32*	8 <sup>†</sup>	37 <sup>†</sup>	999,000	Non-bias representation of templates for genome and seq-based applications	High error rates compared with other reversible terminator chemistries	Seq-based methods	91
Pacific Biosciences (target release: 2010)	Frag only/ single molecule	Real-time	964*	N/A	N/A	N/A	Has the greatest potential for reads exceeding 1 kb	Highest error rates compared with other NGS chemistries	Full-length transcriptome sequencing; complements other resequencing efforts in discovering large structural variants and haplotype blocks	S. Turner, pers. comm.

# Then, by 2014, an ecosystem of options erupted

Table 1: Types of High-Throughput Sequencing Technologies

Optical Sequencing					
Platform	Instrument	Template Preparation	Chemistry	Average Length	Longest Read
Illumina	HiSeq2500	BridgePCR/cluster	Rev. Term., SBS	100	150
Illumina	HiSeq2000	BridgePCR/cluster	Rev. Term., SBS	100	150
Illumina	MiSeq	BridgePCR/cluster	Rev. Term., SBS	250	300
GnuBio	GnuBio	emPCR	Hyb-Assist Sequencing	1000*	64,000*
Life Technologies	SOLID 5500	emPCR	Seq. by Lig.	75	100
LaserGen	LaserGen	emPCR	Rev. Term., SBS	25*	100*
Pacific Biosciences	RS	Polymerase Binding	Real-time	1800	15,000
454	Titanium	emPCR	PyroSequencing	650	1100
454	Junior	emPCR	PyroSequencing	400	650
Helicos	Heliscope	adaptor ligation	Rev. Term., SBS	35	57
Intelligent BioSystems	MAX-Seq	Rolony Amplification	Two-Step SBS (label/unlabel)	2x100	300
Intelligent BioSystems	MINI-20	Rolony Amplification	Two-Step SBS (label/unlabel)	2x100	300
ZS Genetics	N/A	Atomic Labeling	Electron Microscope	N/A	N/A
Halcyon Molecular	N/A	N/A	Direct Observation of DNA	N/A	N/A
Electrical Sequencing					
Platform	Instrument	Template Preparation	Chemistry	Average Length	Longest Read
IBM DNA Transistor	N/A	none	Microchip Nanopore	N/A	N/A
NABsys	N/A	none	Nanochannel	N/A	N/A
Bionanogenomics	N/A	anneal 7mers	Nanochannel	N/A	N/A
Life Technologies	PGM	emPCR	Semi-conductor	150	300
Life Technologies	Proton	emPCR	Semi-conductor	120	240
Life Technologies	Proton 2	emPCR	Semi-conductor	400*	800*
Genia	N/A	none	Protein Nanopore (a-hemalysin)	N/A	N/A
Oxford Nanopore	MinION	none	Protein Nanopore	10,000	10,000*
Oxford Nanopore	GridION 2K	none	Protein Nanopore	10,000	500,000*
Oxford Nanopore	GridION 8K	none	Protein Nanopore	10,000	500,000*

\*Values are estimates from companies that have not yet released actual data

# Coming of age: ten years of next-generation sequencing technologies

---

*Sara Goodwin<sup>1</sup>, John D. McPherson<sup>2</sup> and W. Richard McCombie<sup>1</sup>*

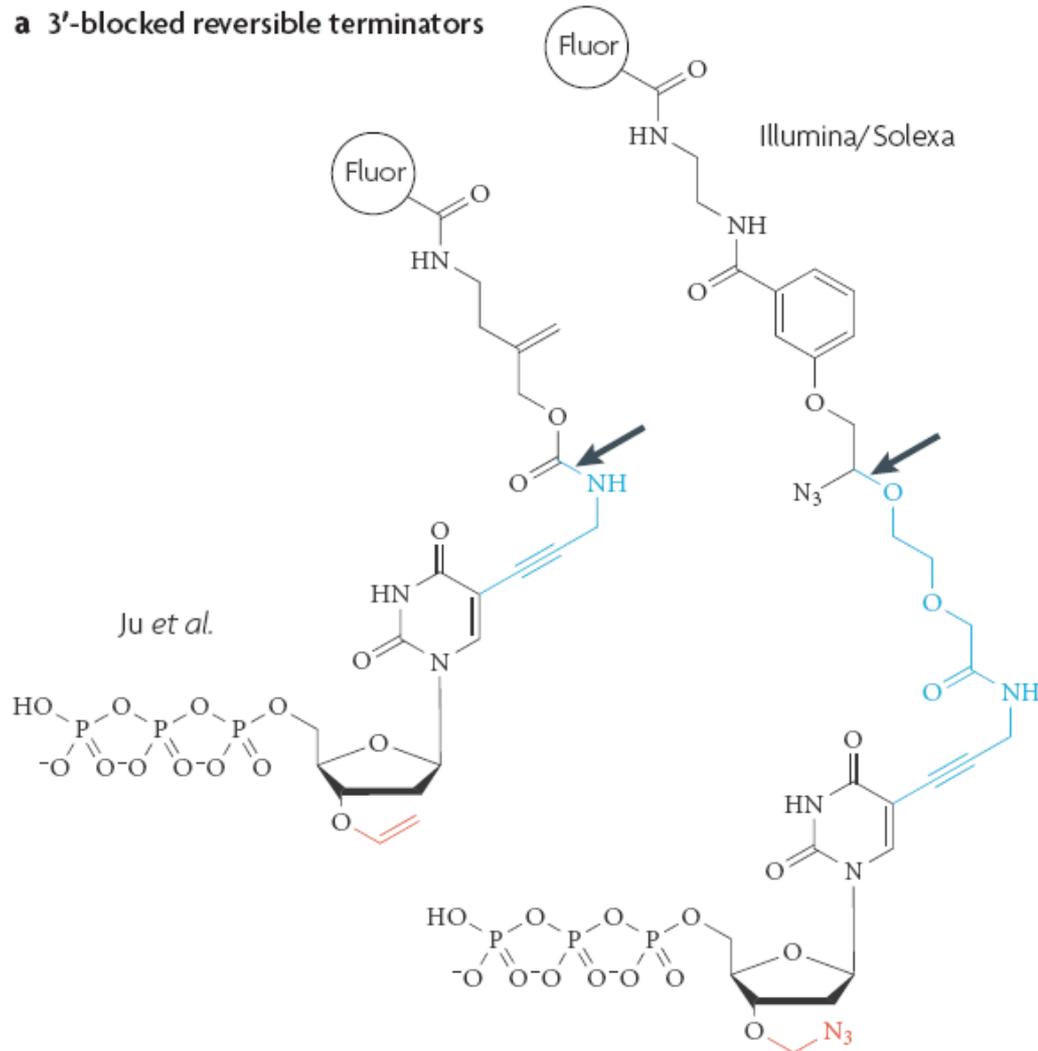
**Abstract** | Since the completion of the human genome project in 2003, extraordinary progress has been made in genome sequencing technologies, which has led to a decreased cost per megabase and an increase in the number and diversity of sequenced genomes. An astonishing complexity of genome architecture has been revealed, bringing these sequencing technologies to even greater advancements. Some approaches maximize the number of bases sequenced in the least amount of time, generating a wealth of data that can be used to understand increasingly complex phenotypes. Alternatively, other approaches now aim to sequence longer contiguous pieces of DNA, which are essential for resolving structurally complex regions. These and other strategies are providing researchers and clinicians a variety of tools to probe genomes in greater depth, leading to an enhanced understanding of how genome sequence variants underlie phenotype and disease.



Consideration of each platform

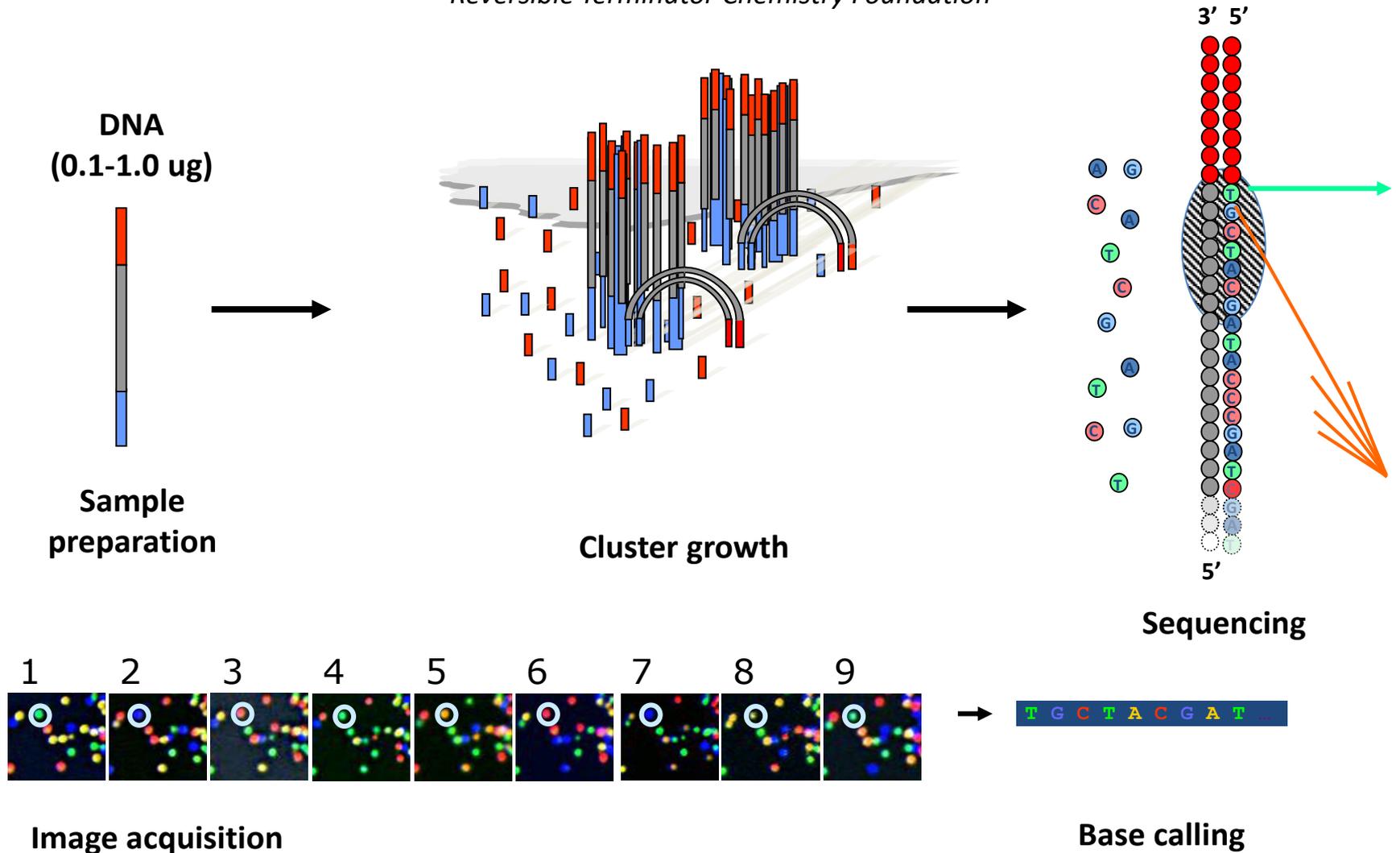
## 2. Reversible Terminator Bases are Essential Technology Used in Many Chemistries

a 3'-blocked reversible terminators



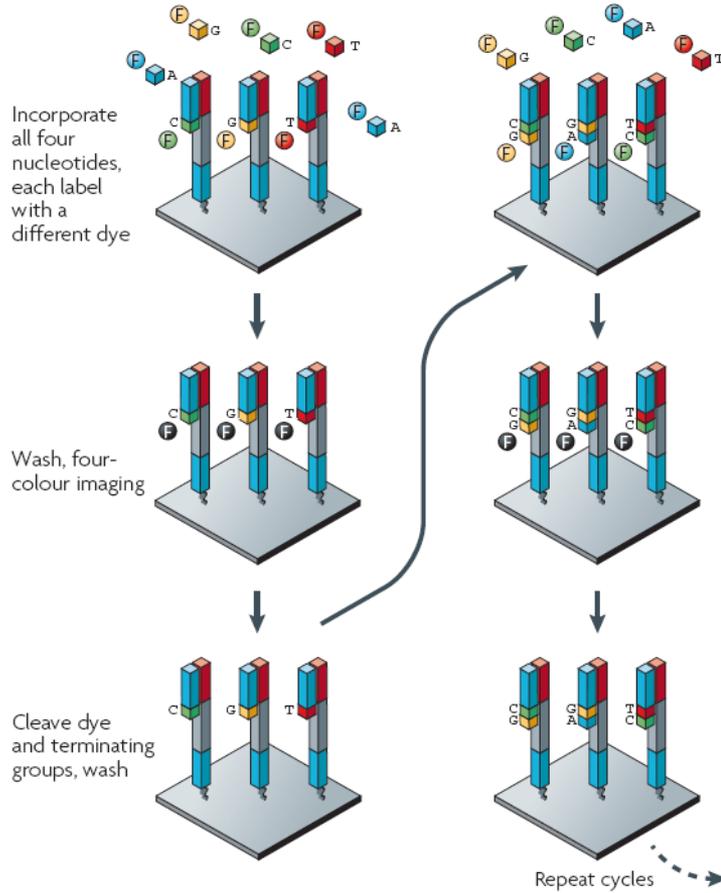
# Illumina SBS Technology

*Reversible Terminator Chemistry Foundation*

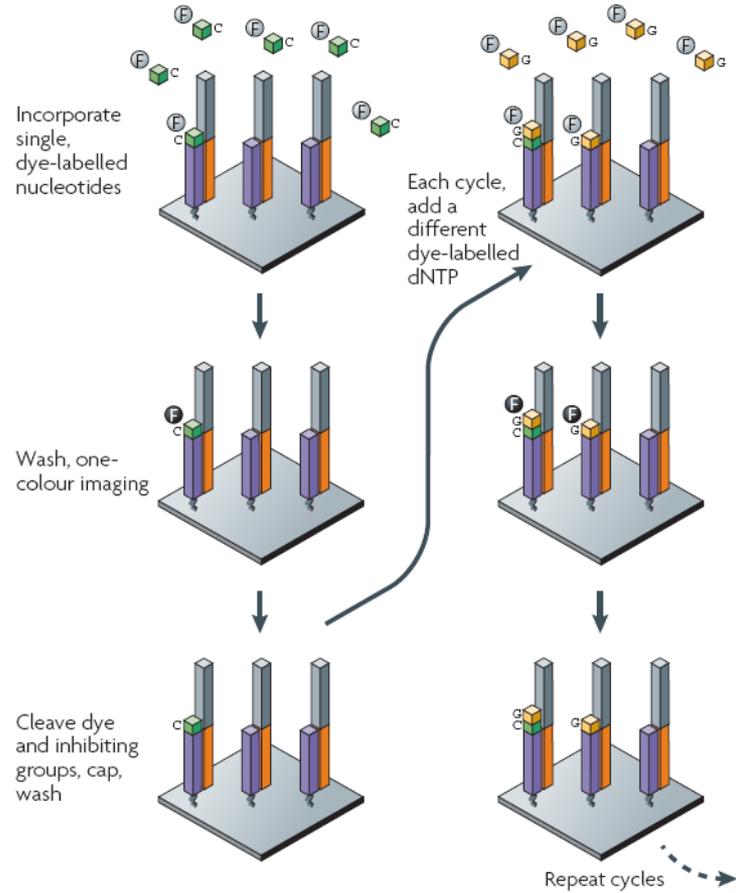


# Sequencing by Synthesis (SBS)

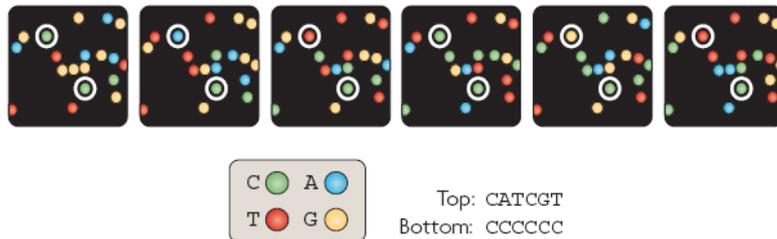
**a** Illumina/Solexa — Reversible terminators



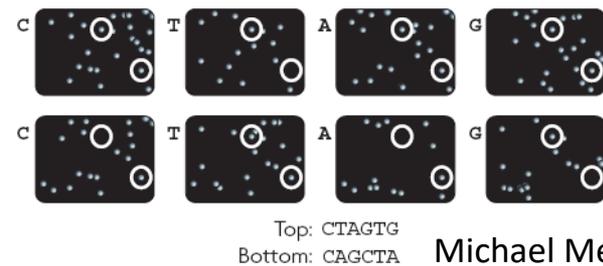
**c** Helicos BioSciences — Reversible terminators



**b**



**d**



# Now three kinds of chemistry

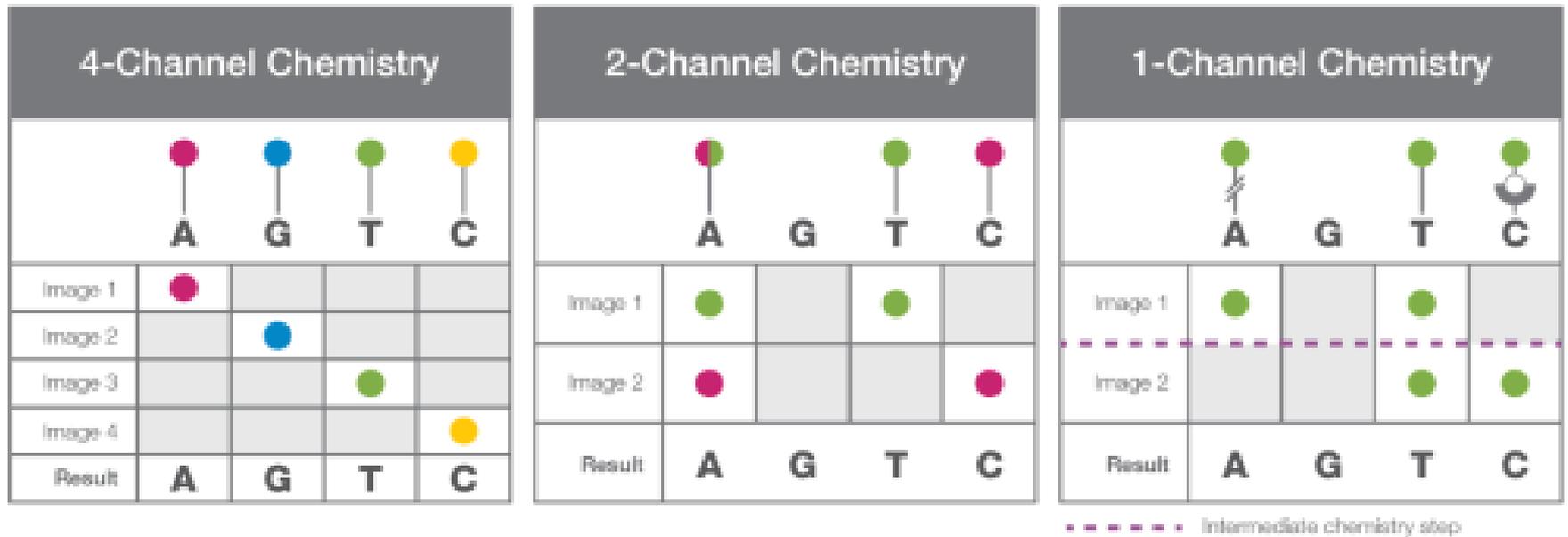
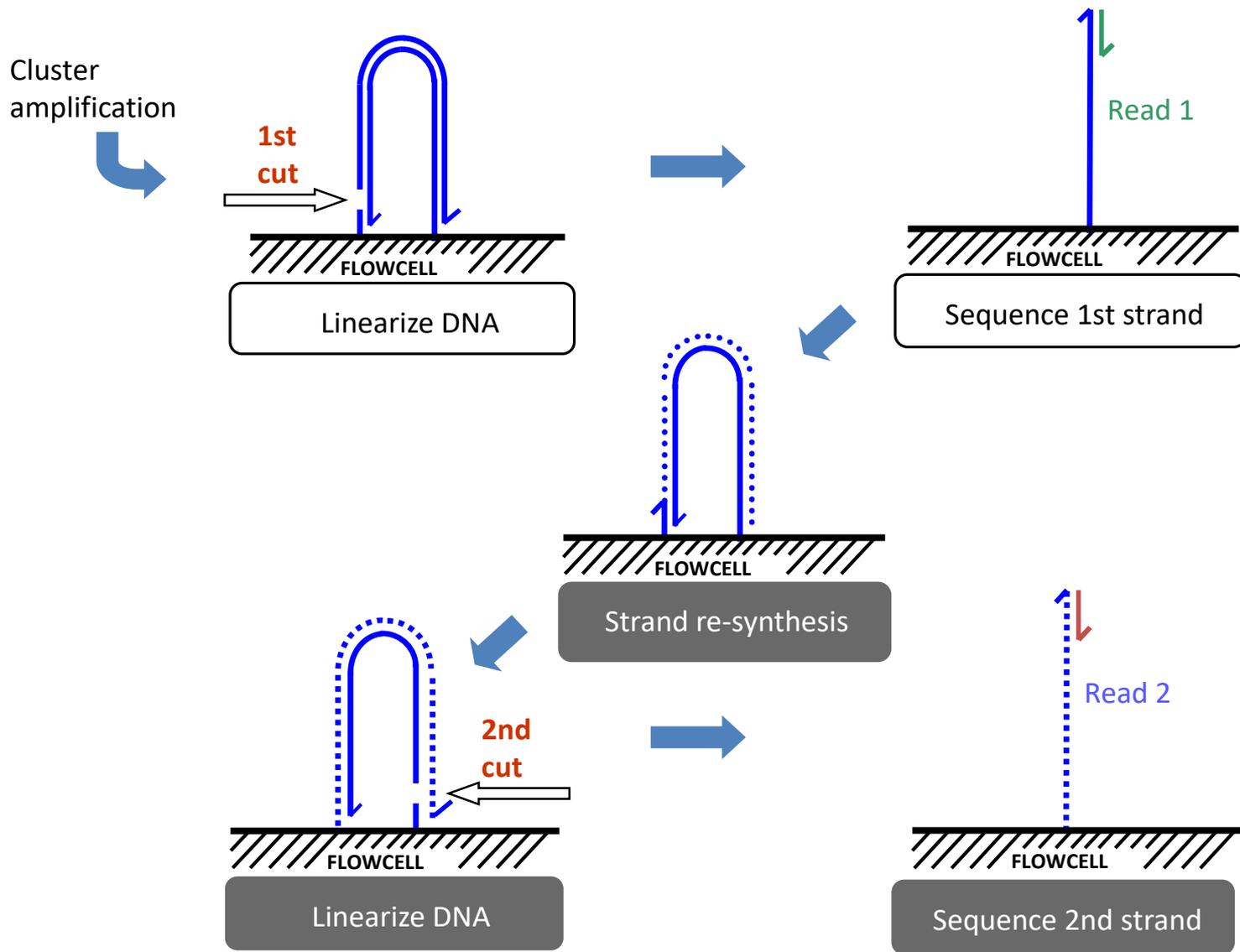
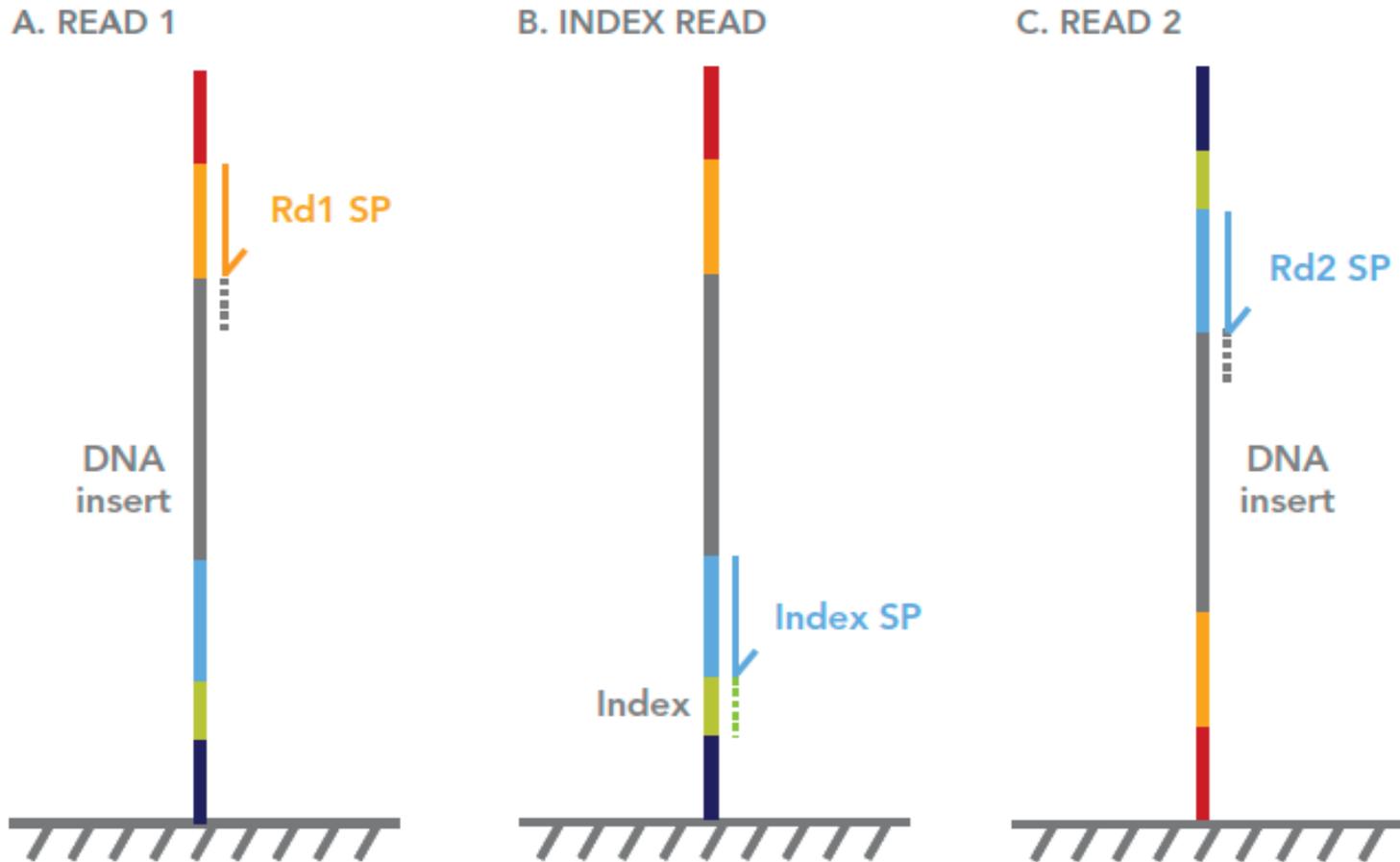


Figure 2: Four-, Two-, and One-Channel Chemistry—Four-channel chemistry uses a mixture of nucleotides labeled with four different fluorescent dyes. Two-channel chemistry uses two different fluorescent dyes, and one-channel chemistry uses only one dye. The images are processed by image analysis software to determine nucleotide identity.

# Paired-End Sequencing allows for two looks at a sequence



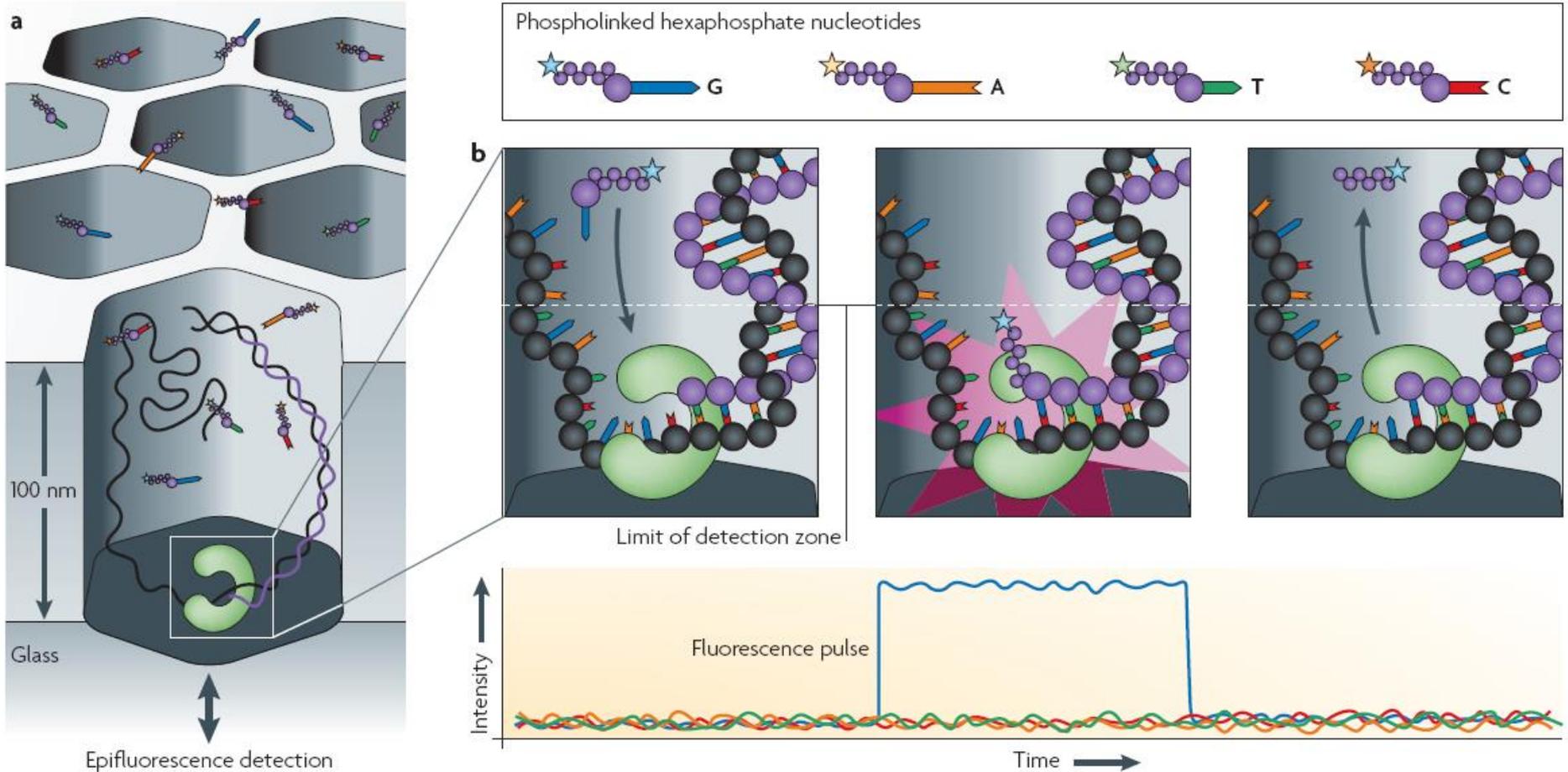
# Indexed sequencing method is now standard for single and paired reads



# Pacific Biosciences

## Single Molecule Real-Time (SMRT) Sequencing

Pacific Biosciences — Real-time sequencing

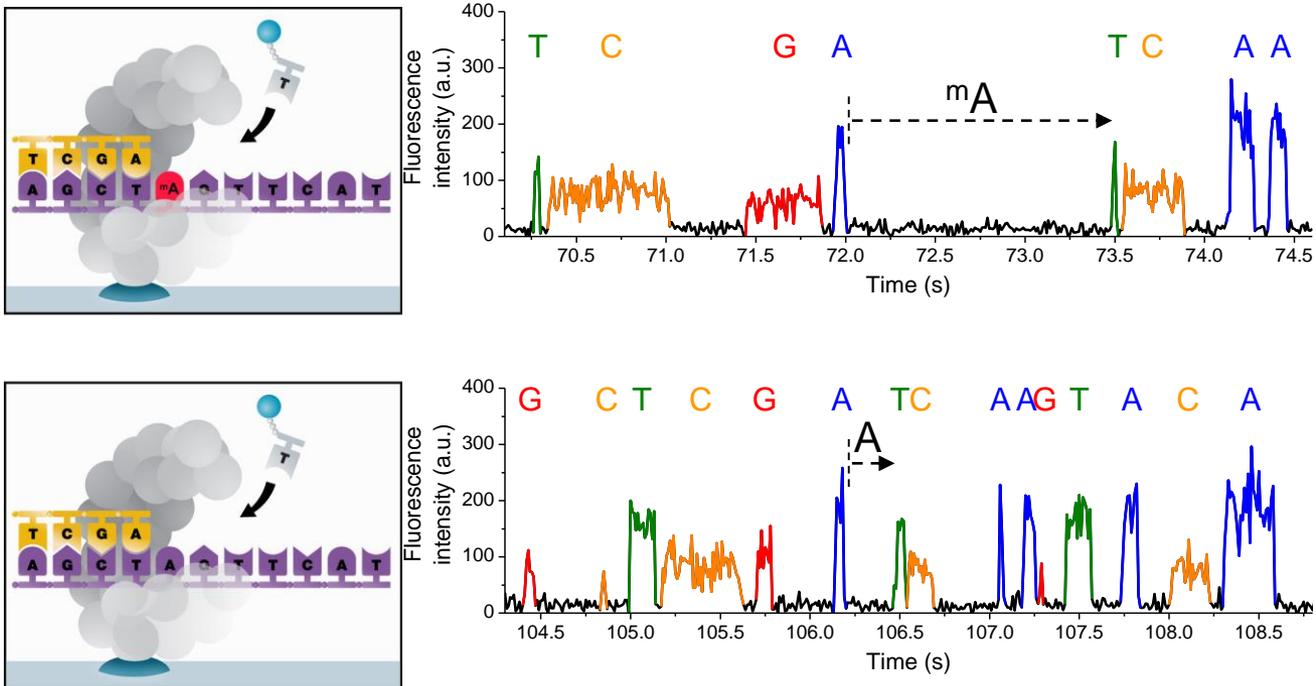




# Single Molecule Kinetics Allow for the Direct Detection of Methylation

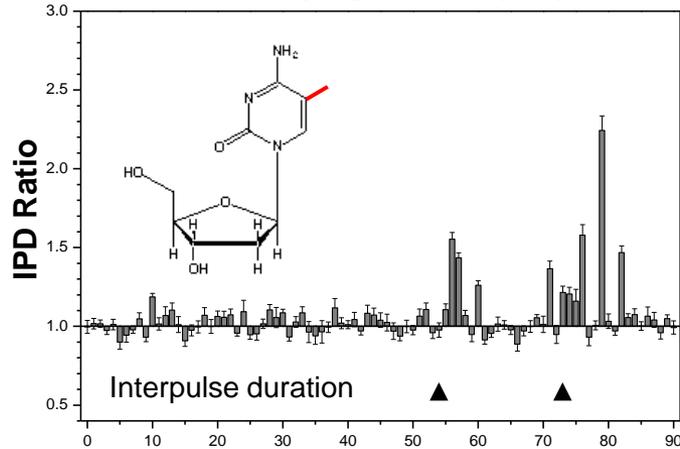
Approach: Kinetic detection of methylated bases during SMRT DNA sequencing

Example: N<sup>6</sup>-methyladenosine (m<sup>A</sup>)

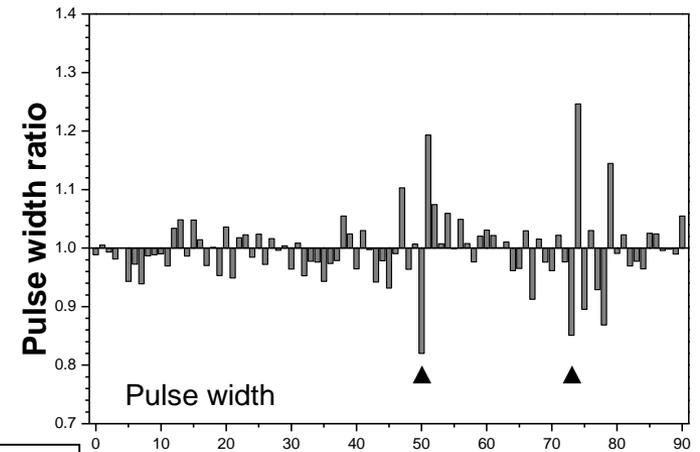
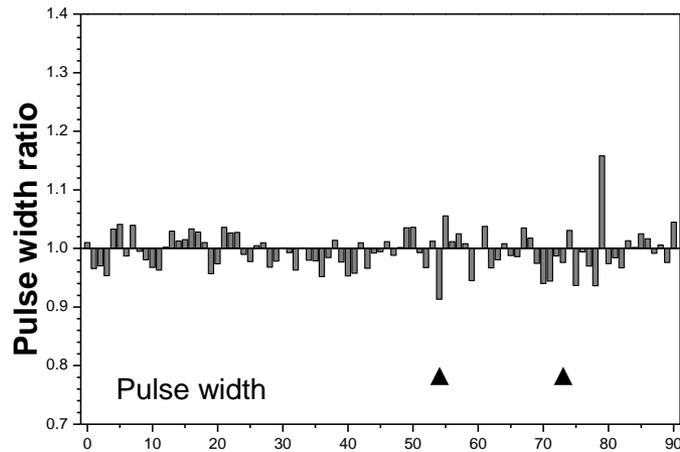
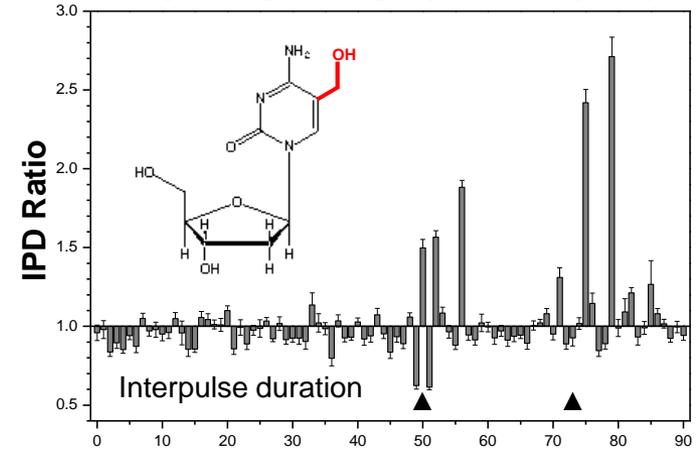


# Kinetics can detect other base modifications

## 5-methylcytosine (mC)



## 5-hydroxymethylcytosine (hmC)

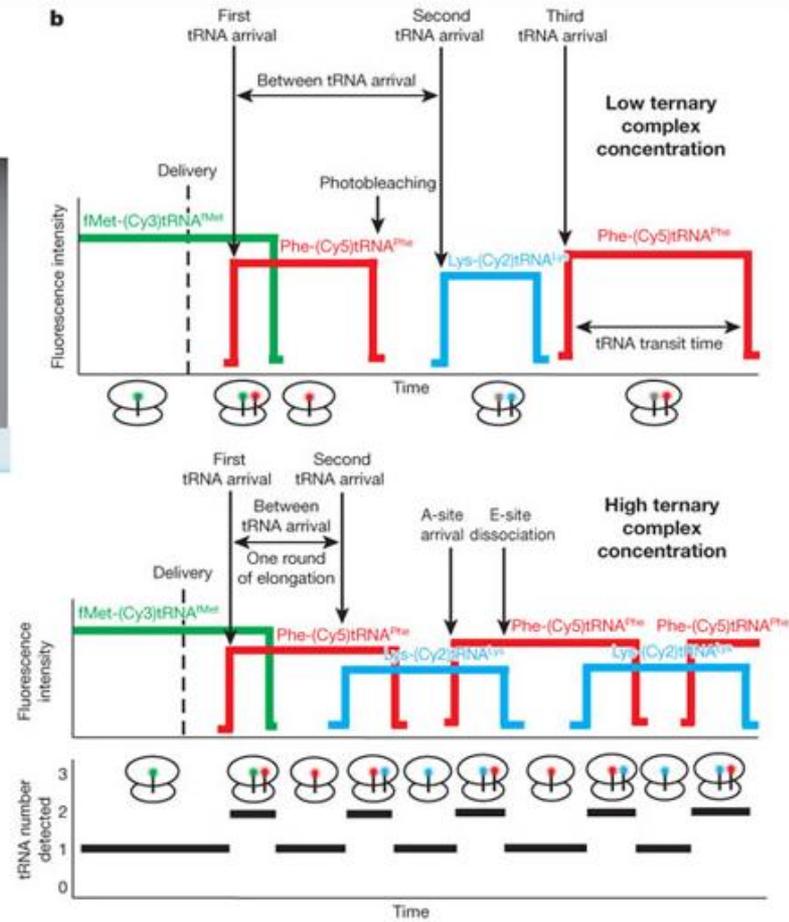
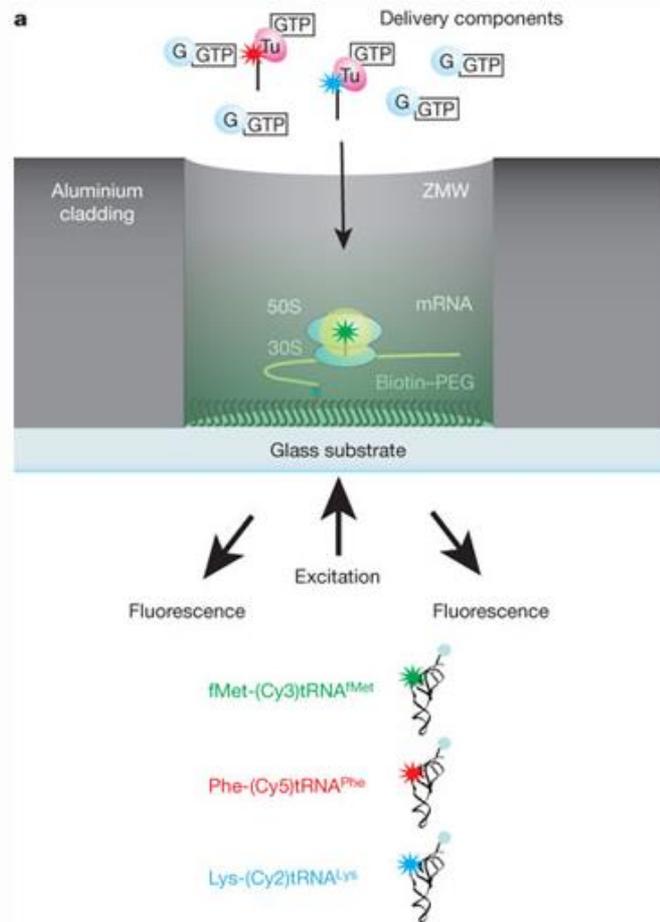


▲ = Methylated position

DNA Template Position

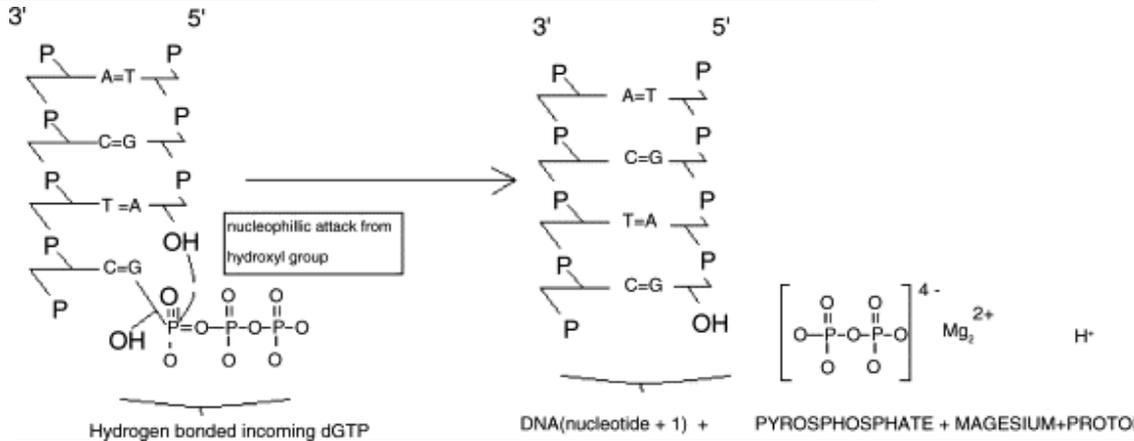
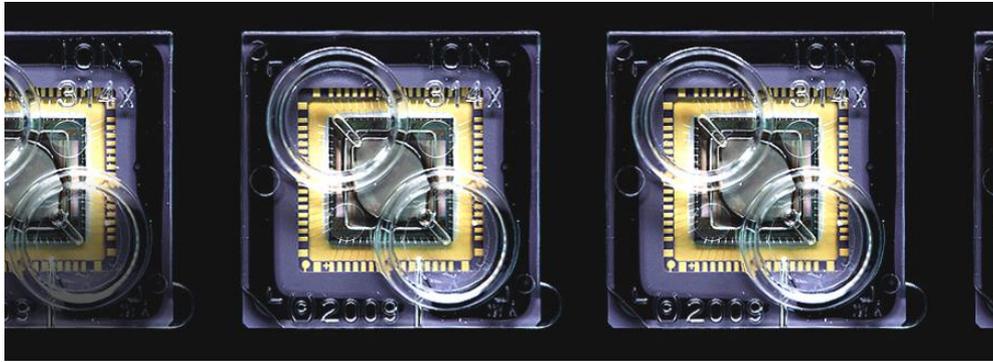
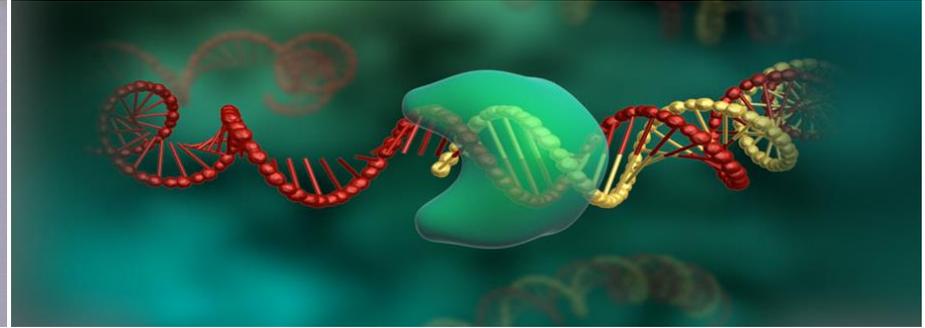
DNA Template Position

# Kinetics allow one to watch protein translation as it occurs



# “Post-Light,” Semi-Conductor Sequencing:

Thermo Fisher’s Personal Genome Machine (PGM), the Proton I and Proton II, and S5



Essentially,  
Millions of  
very small  
pH meters

Purushothaman *et al*, 2005  
IonTorrent, Inc.

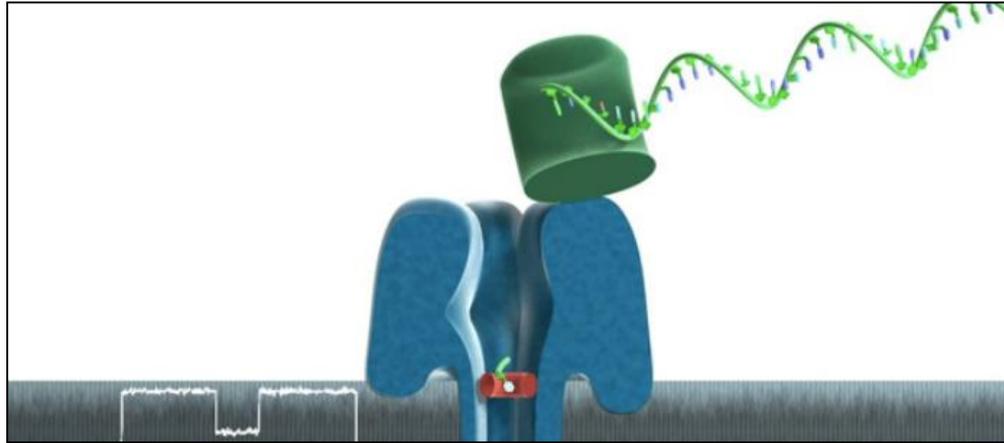
# Latest Ion Platforms

## Thermo Fisher's Ion S5 & S5 XL

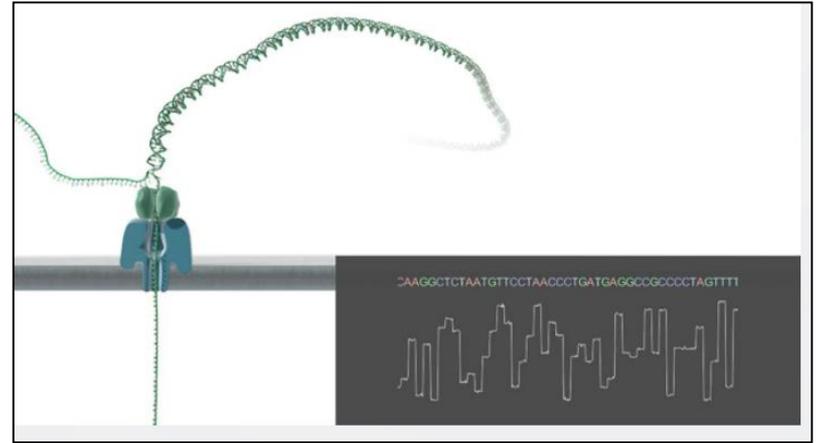




# 2014: Sequencing with a protein nanopore



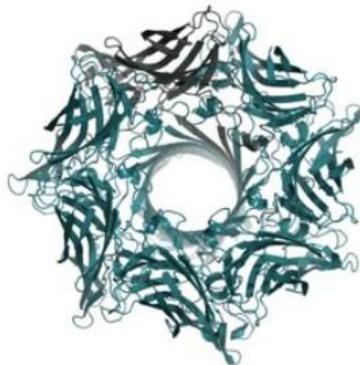
Exonuclease-Seq



Strand-Seq



MinION

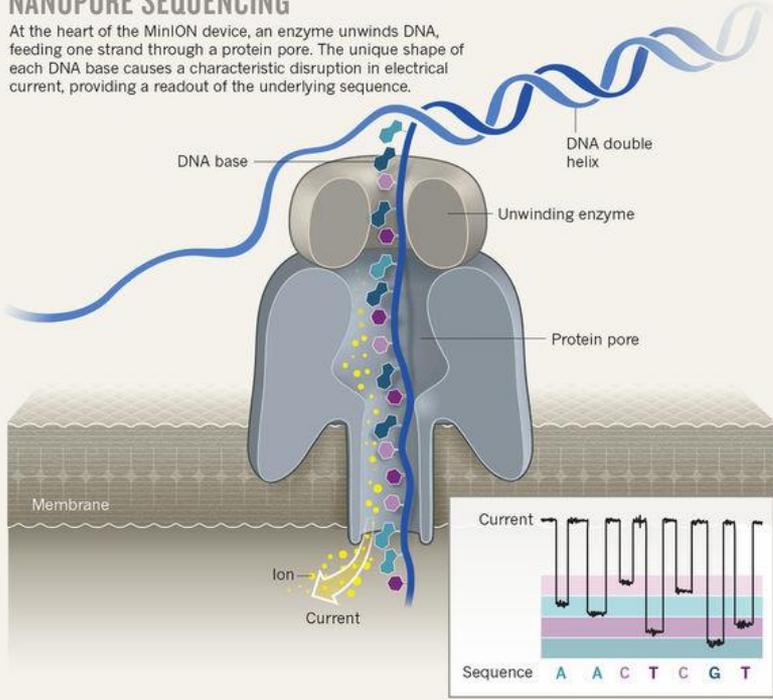


PromethION

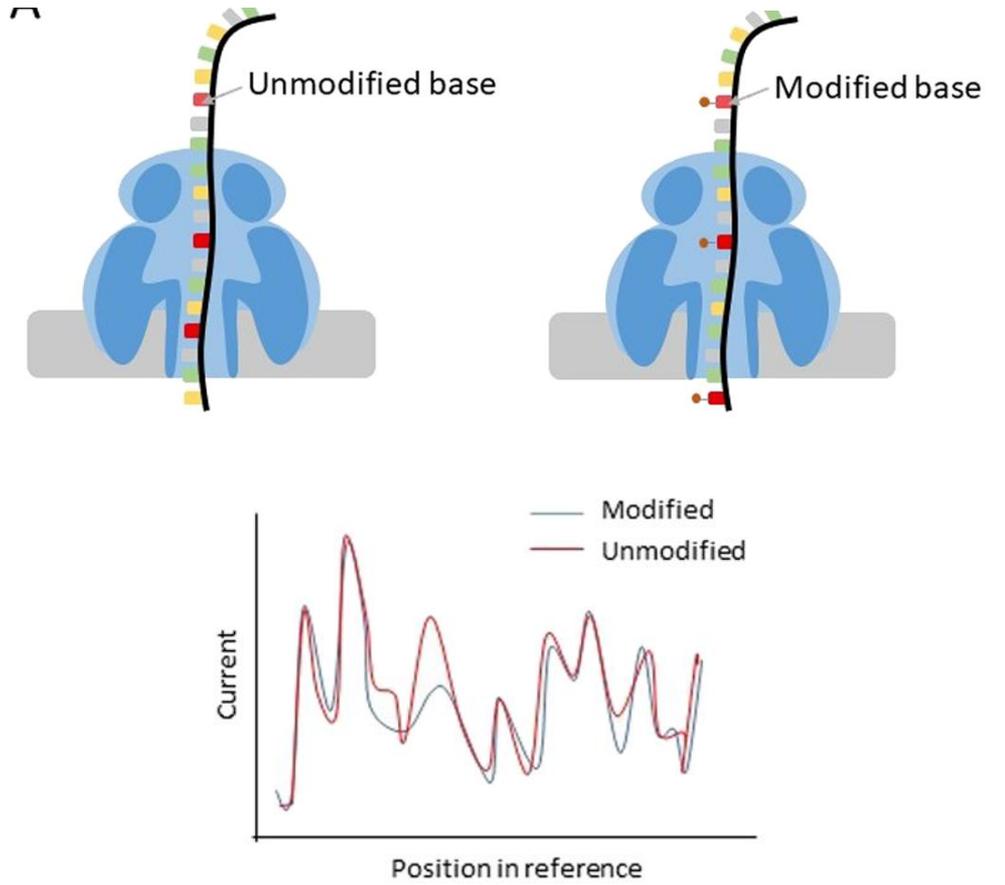


## NANOPORE SEQUENCING

At the heart of the MinION device, an enzyme unwinds DNA, feeding one strand through a protein pore. The unique shape of each DNA base causes a characteristic disruption in electrical current, providing a readout of the underlying sequence.



<http://blogs.nature.com/naturejobs/2017/10/16/techblog-the-nanopore-toolbox/>



# 2022

## Products & Services



### Sequencing platforms

Learn more

### Consumables

Flow cells



Kits & sample prep



### Research

Real-time DNA and RNA sequencing — from portable to high-throughput devices.



### IVD testing

LamPORE — rapid, low-cost, highly scalable detection of SARS-CoV-2.



### Q-Line

Locked-down, research-validated devices for applied sequencing applications.

<https://nanoporetech.com/>

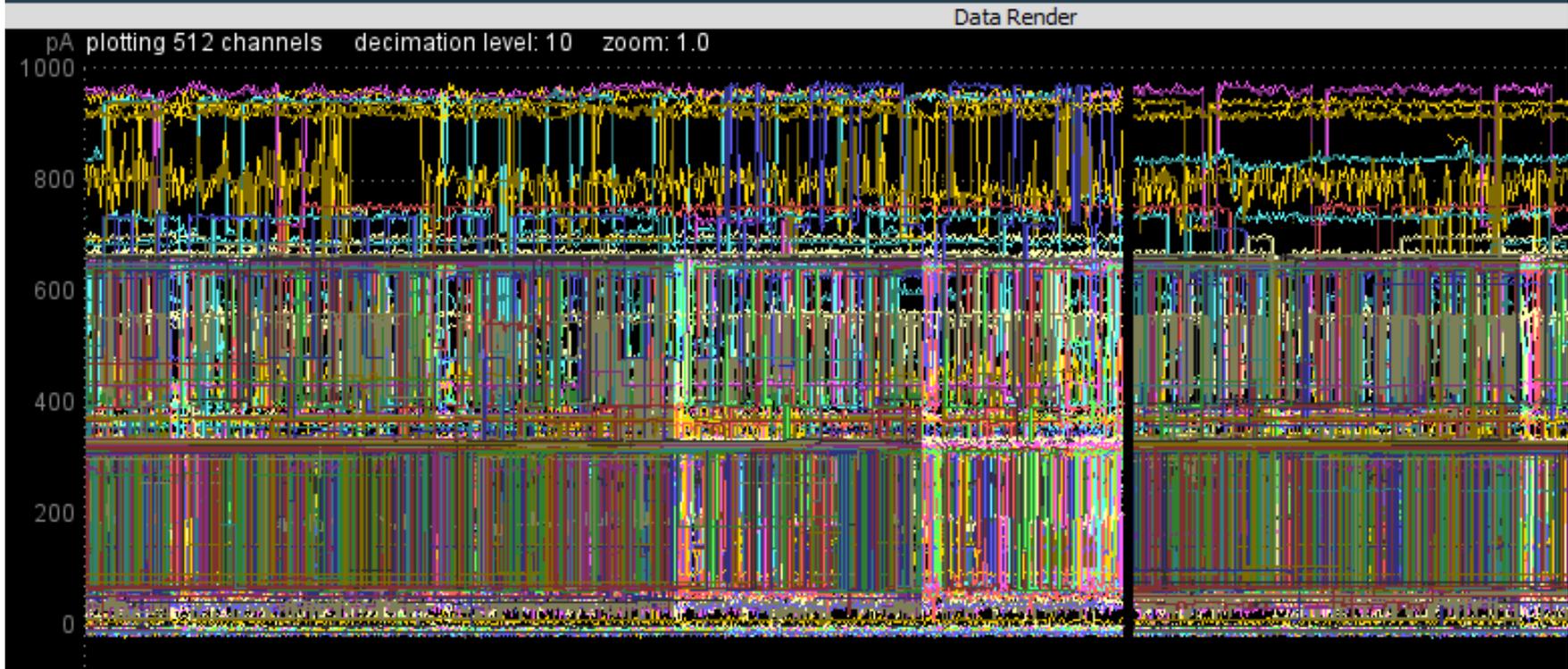


# They are small



# Base space is now “squiggle space”

Status	●	Exp. Time	478s	Asic Status	●				
Acquisition	●	Yield	4710240	Asic	37.1°C	Bias Voltage	-180mV		
Analysis	●	Sample ID	BN_011_34887GT_wu598_new	Heatsink	36.5°C	MinION ID	MN02301	Asic ID	37299
Write	●	Data Set	NRCHBS-WDL31403_BN_011_34887GT_w...						



# You can do it anywhere



**nature**  
International journal of science

Letter | Published: 03 February 2016

## Real-time, portable genome sequencing for Ebola surveillance

Joshua Quick, Nicholas J. Loman  [...] Miles W. Carroll

*Nature* **530**, 228–232 (11 February 2016) | [Download Citation](#) ↓

<https://www.nature.com/articles/nature16996>



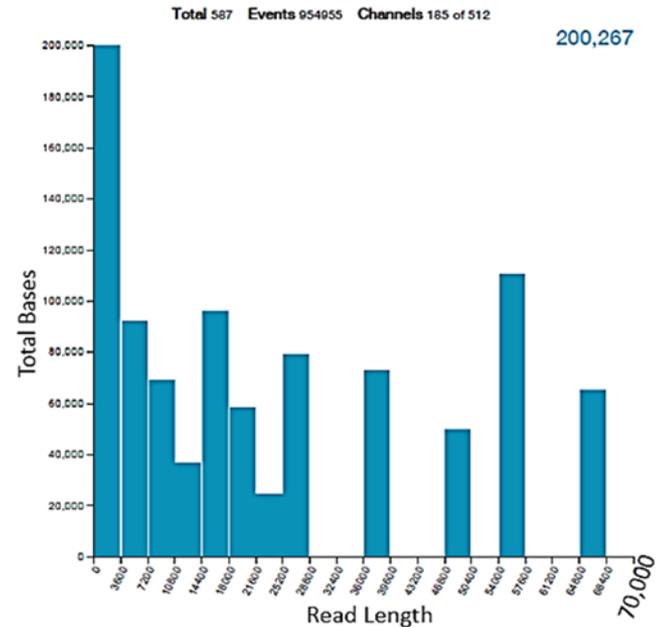
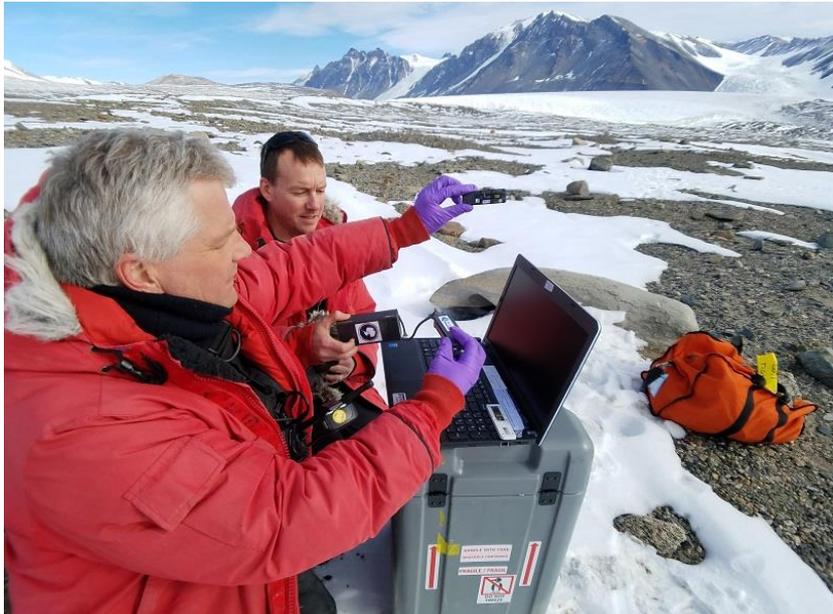


Scott Tighe

# Lake Fryxell, Antarctica

## Scott Tighe

Sequencing HW DNA in the field with the Oxford Nanopore  
Sarah Johnson (PI) expedition G062 team



J. Biomol. Tech. 2017 Apr; 18(17):2801-2809  
Published online 2017 Mar 22 doi: 10.1039/c6jbt00150g

PMCID: PMC5362188

### Real-Time DNA Sequencing in the Antarctic Dry Valleys Using the Oxford Nanopore Sequencer

Sarah S. Johnson,<sup>1,2</sup> Elena Zaikova,<sup>1</sup> David S. Goerlitz,<sup>3</sup> Yu Bai,<sup>1</sup> and Scott W. Tighe<sup>4</sup>

Author information ► Copyright and License information ►

Abstract

### ARTICLE

#### Genomic Methods and Microbiological Technologies for Profiling Novel and Extreme Environments for the Extreme Microbiome Project (XMP)

Scott Tighe,<sup>1,2,4</sup> Ebrahim Afshinneko,<sup>2,5,6</sup> Tara M. Rock,<sup>5</sup> Ken McGrath,<sup>6</sup> Noah Alexander,<sup>2,3</sup> Alexa McIntyre,<sup>2,3</sup> Sofia Abanuddin,<sup>2,3</sup> Daniela Bezdán,<sup>2,3</sup> Stefan J. Green,<sup>7</sup> Samantha Jey,<sup>8</sup> Sarah Stewart Johnson,<sup>9</sup> Don A. Baldwin,<sup>10</sup> Nathan Bivens,<sup>11</sup> Nadim Ajami,<sup>12,13</sup> Joseph R. Carmical,<sup>12,13</sup> Ian Charold Herriott,<sup>14</sup> Rita Colwell,<sup>15</sup> Mohamed Donia,<sup>16</sup> Jonathan Fox,<sup>2,5,17</sup> Nick Greenfield,<sup>18</sup> Tim Hunter,<sup>1</sup> Jessica Hoffman,<sup>1</sup> Joshua Hyman,<sup>17</sup> Ellen Jorgensen,<sup>20</sup> Diana Krauczyk,<sup>21</sup> Jodie Lee,<sup>22</sup> Shawn Levy,<sup>23</sup> Natalia Garcia-Reyero,<sup>24</sup> Matthew Settle,<sup>25</sup> Kelley Thomas,<sup>26</sup> Felipe Gómez,<sup>27</sup> Lynn Schriml,<sup>28,29</sup> Nikos Kyrpides,<sup>30</sup> Elena Zaikova,<sup>1</sup> Jon Penterman,<sup>31</sup> and Christopher E. Mason<sup>2,3,32</sup>

# Zero-G Pipetting: Hardest Lab Job Ever



Dr. Andrew Feinberg

NATURE | NEWS



## Zero-gravity genomics passes first test

Two experiments demonstrate sample transfer and sequencing in a low-gravity environment.

[Chris Cesare](#)

13 October 2015

 [Rights & Permissions](#)

After 160 swoops in NASA's zero-gravity aeroplane, researchers have the first evidence that genetic sequencing can be done in space.





MENU ▾

npj | Microgravity



Search



E-alert



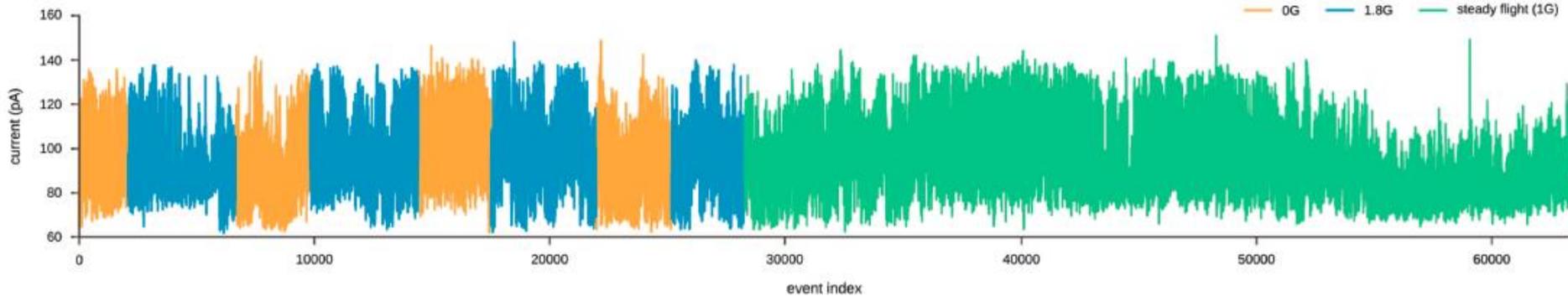
Submit



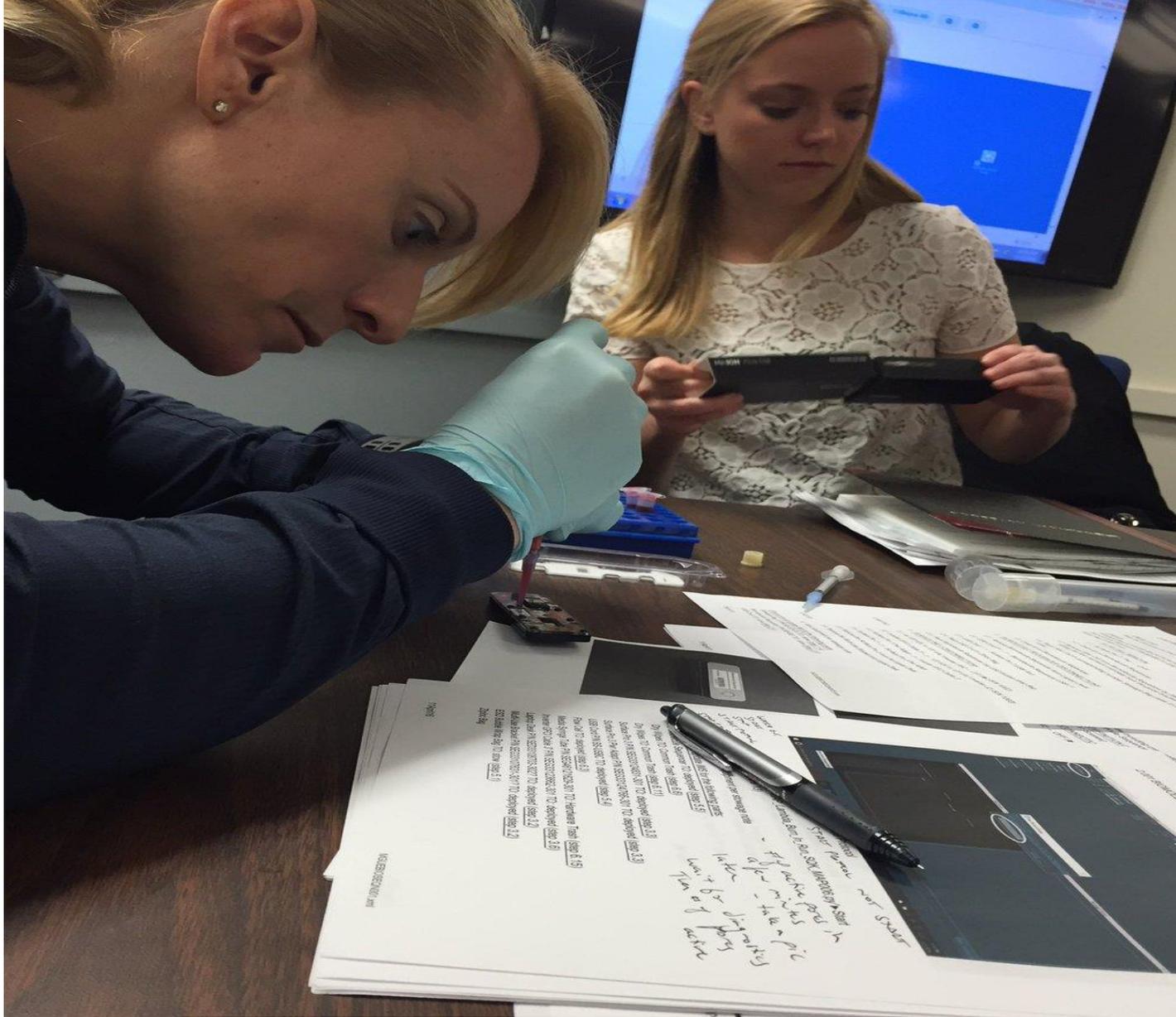
Login

# DNA sequencing in space: Nanopores ready for liftoff

Results from the first DNA sequencing experiments performed in microgravity reveal a promising future for portable 'nanopore' devices in space missions. Read the paper in full.



McIntyre ABR et al., *Nature Microgravity*, 2016.



**Christopher Mason** @mason\_lab ·

Preparing for sequencing in space! @Astro\_Kate7 @NASA

@ScientistAaronB 450uL in one load should work w/ @nanopore

# SpaceX CRS-7 blows up



National Aeronautics and Space Administration

Office of the Administrator  
Washington, DC 20546-0001



Dr. Christopher Mason  
Weill Cornell Medical College  
1300 York Ave.  
New York, NY 10065

Dear Dr. Mason:

As NASA astronaut Scott Kelley tweeted on Sunday, June 28, 2015, "space is hard."

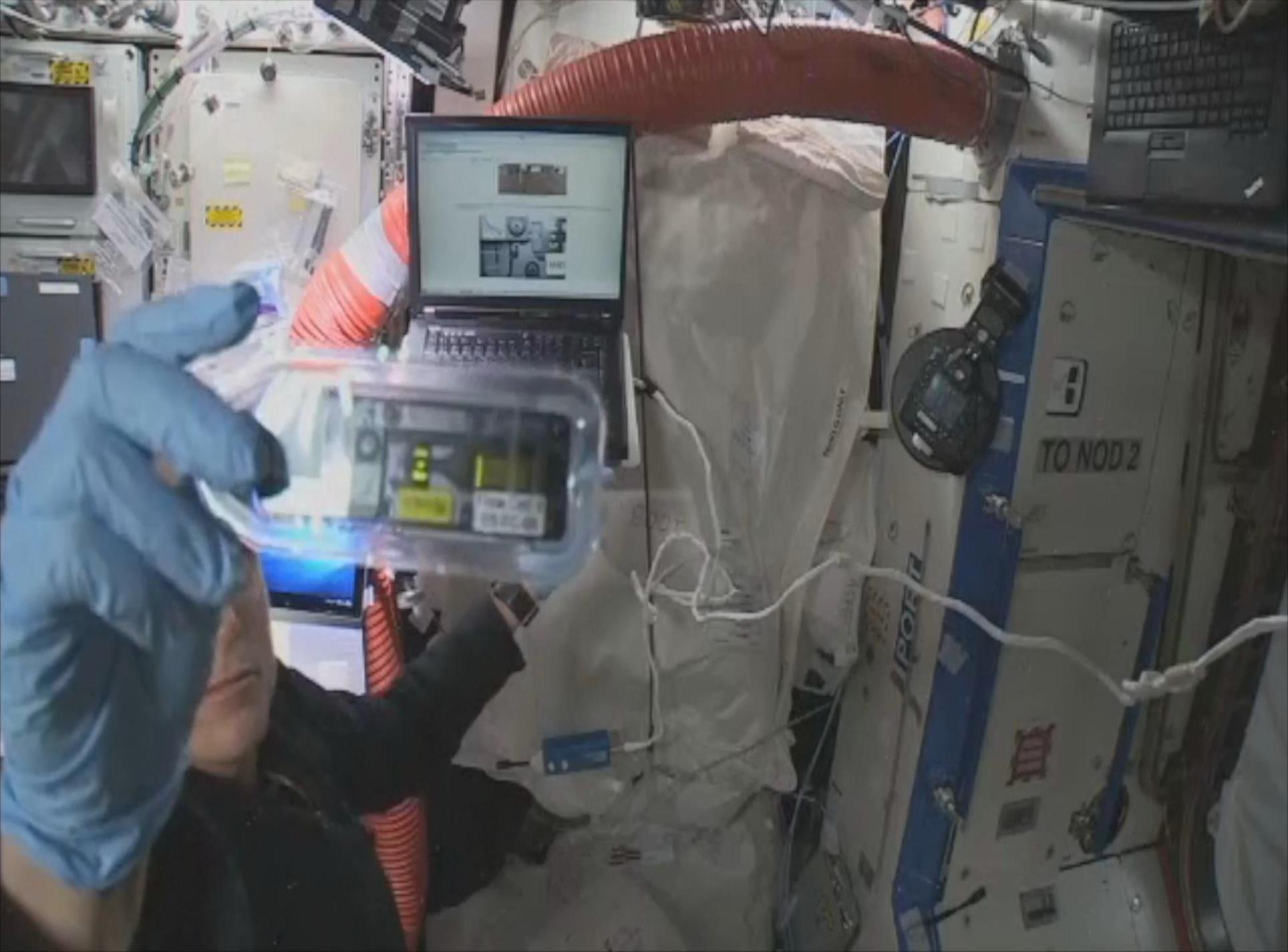
Speaking as a fellow researcher, I can only imagine how devastated you must be feeling right now with the loss of SpaceX's CRS-7. I am saddened and disappointed too. I am sure that the tremendous honor of being selected to have your experiment flown on the International Space Station is of little solace after the loss of months, and perhaps even years, of hard work.

I am writing to encourage you – and in fact, to urge you – to continue your inquiry. The story of space exploration is the story of people just like you who meet adversity, head on, with determination and scientific and technological advancement. If you think about it, virtually every major innovation and technological breakthrough in human history has been the product of many different stops and starts; learning and being better because of failures and setbacks and, ultimately, enhanced knowledge and moving forward.



SpaceX CRS-9: perfect launch  
and booster return  
July 18, 2016







## Latest

## Related



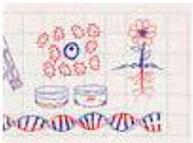
Weekly Recap From the Expedition Lead Scientist

*5 days ago*



Weekly Recap From the Expedition Lead Scientist

*13 days ago*



Biological Sciences on the International Space Station

*19 days ago*



SAGE III to Look Back at Earth's Atmospheric 'Sunscreen'

*19 days ago*



Weekly Recap From the Expedition Lead Scientist

*19 days ago*



Weekly Recap From the Expedition Lead Scientist

*24 days ago*



Weekly Recap From the Expedition Lead Scientist

*a month ago*

## Space Station



Aug. 29, 2016

# First DNA Sequencing in Space a Game Changer

For the first time ever, DNA was successfully sequenced in microgravity as part of the [Biomolecule Sequencer](#) experiment performed by NASA astronaut Kate Rubins this weekend aboard the [International Space Station](#). The ability to sequence the DNA of living organisms in space opens a whole new world of scientific and medical possibilities. Scientists consider it a game changer.

DNA, or deoxyribonucleic acid, contains the instructions each cell in an organism on Earth needs to live. These instructions are represented by the letters A, G, C and T, which stand for the four chemical bases of DNA, adenine, guanine, cytosine, and thymine. Both the number and arrangement of these bases differ among organisms, so their order, or sequence, can be used to identify a specific organism.





spasmunkey

@spasmunkey



Following

Great to see this team at work from training to operations at "the dawn of genomics...in space"  
#AstroKate



RETWEETS

4

LIKES

12



9:40 PM - 29 Aug 2016

Houston, TX

You, Aaron Burton, Kristen John and 3 others



From zero to one billion: sequencing the one billionth base pair of DNA in space. [go.nasa.gov/2bV2UnD](https://go.nasa.gov/2bV2UnD)



**sequencing the one billionth base pair of DNA**

Clip from NASA TV

RETWEETS

**123**

LIKES

**185**

Bus Lon Dor Elai Alfc Oliv Jee  Lita

3:28 PM - 14 Sep 2016

flight

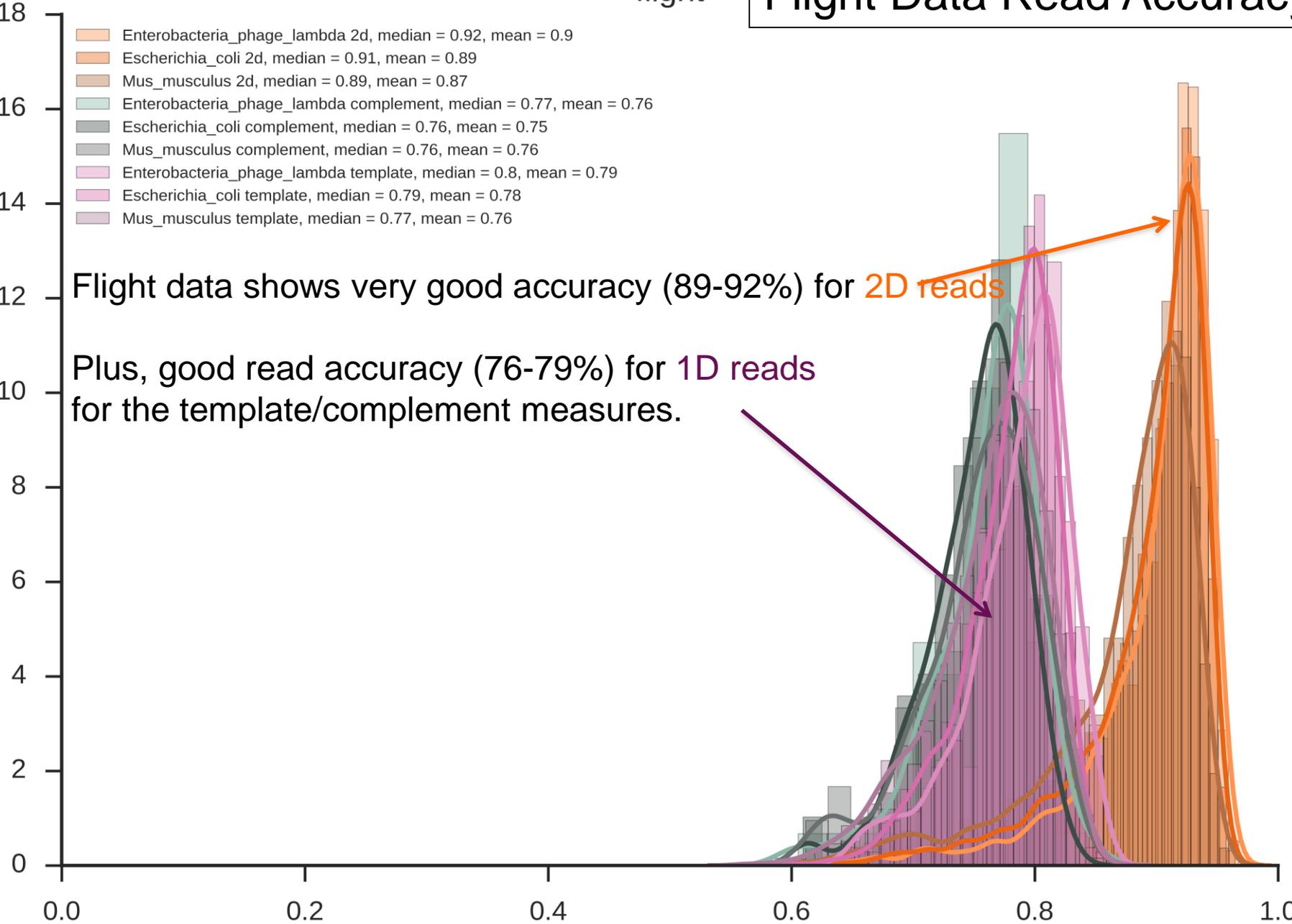
# Flight Data Read Accuracy

- Enterobacteria\_phage\_lambda 2d, median = 0.92, mean = 0.9
- Escherichia\_coli 2d, median = 0.91, mean = 0.89
- Mus\_musculus 2d, median = 0.89, mean = 0.87
- Enterobacteria\_phage\_lambda complement, median = 0.77, mean = 0.76
- Escherichia\_coli complement, median = 0.76, mean = 0.75
- Mus\_musculus complement, median = 0.76, mean = 0.76
- Enterobacteria\_phage\_lambda template, median = 0.8, mean = 0.79
- Escherichia\_coli template, median = 0.79, mean = 0.78
- Mus\_musculus template, median = 0.77, mean = 0.76

(% of reads)

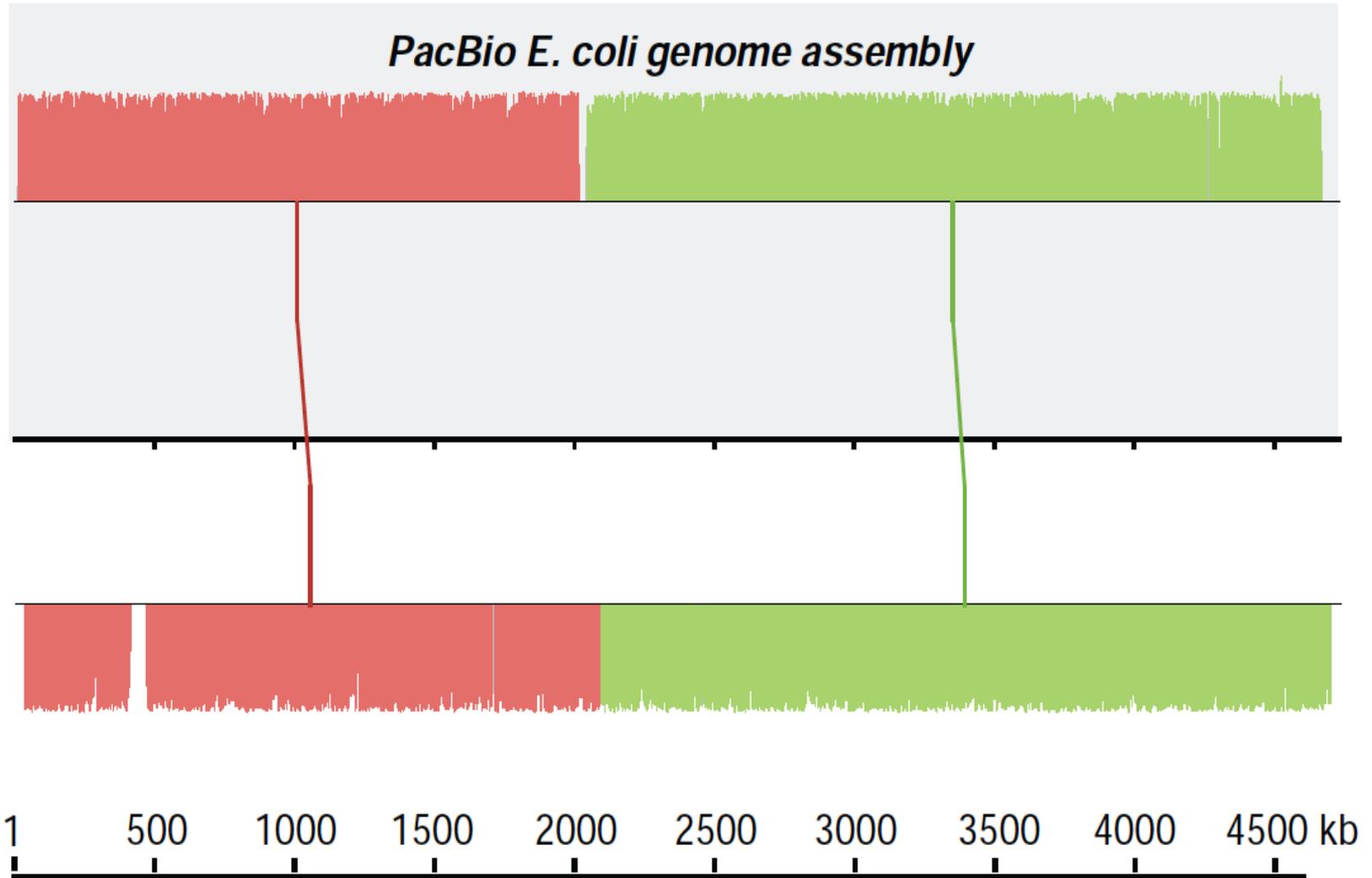
Flight data shows very good accuracy (89-92%) for 2D reads

Plus, good read accuracy (76-79%) for 1D reads for the template/complement measures.



1-2% better than ground data

# Almost perfect when compared to PacBio



# The first genome sequence and assembly off Earth



 Altmetric: 171

[More detail >>](#)

Article | [OPEN](#)

## Nanopore DNA Sequencing and Genome Assembly on the International Space Station

Sarah L. Castro-Wallace, Charles Y. Chiu, Kristen K. John, Sarah E. Stahl, Kathleen H. Rubins, Alexa B. R. McIntyre, Jason P. Dworkin, Mark L. Lupisella, David J. Smith, Douglas J. Botkin, Timothy A. Stephenson, Sissel Juul, Daniel J. Turner, Fernando Izquierdo, Scot Federman, Doug Stryke, Sneha Somasekar, Noah Alexander, Guixia Yu, Christopher E. Mason & Aaron S. Burton 

*Scientific Reports* **7**, Article number: 18022  
(2017)  
doi:10.1038/s41598-017-18364-0

Received: 01 August 2017  
Accepted: 11 December 2017  
Published online: 21 December 2017

<https://www.nature.com/articles/s41598-017-18364-0>

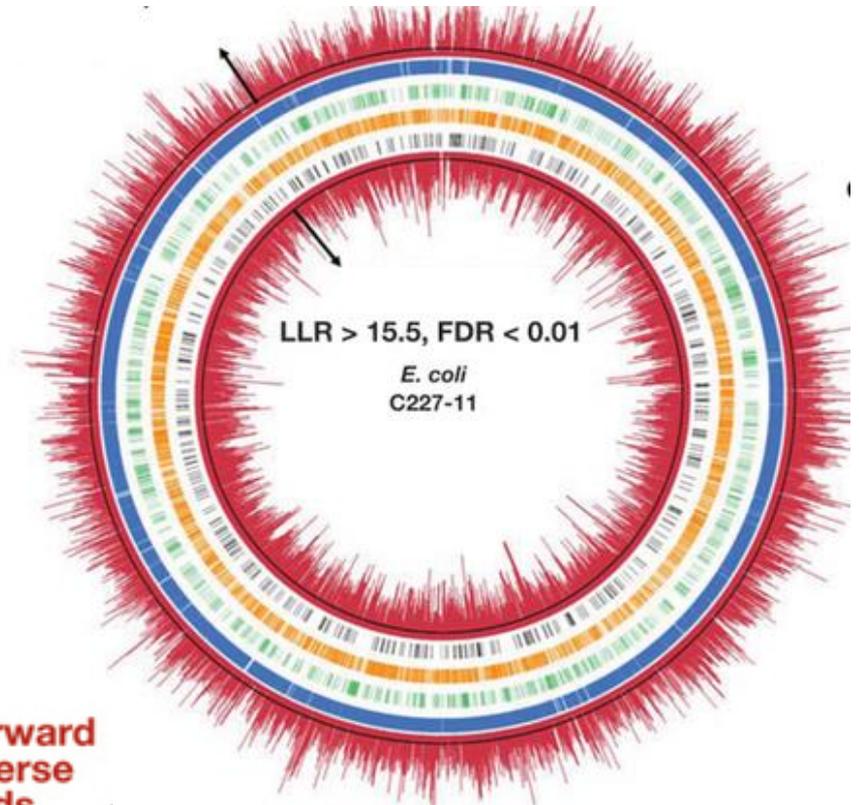
# Bacteria are splattered with epigenetic marks

Genome-wide mapping of methylated adenine residues in pathogenic *Escherichia coli* using single-molecule real-time sequencing

Gang Fang, Diana Munera, David I Friedman, Anjali Mandlik, Michael C Chao, Onureena Banerjee, Zhixing Feng, Bojan Losic, Milind C Mahajan, Omar J Jabado, Gintaras Deikus, Tyson A Clark, Khai Luong, Iain A Murray, Brigid M Davis, Alona Keren-Paz, Andrew Chess, Richard J Roberts, Jonas Koriach, Steve W Turner, Vipin Kumar, Matthew K Waldor & Eric E Schadt

[Affiliations](#) | [Contributions](#) | [Corresponding authors](#)

*Nature Biotechnology* 30, 1232–1239 (2012) | doi:10.1038/nbt.2432



LLRs, forward  
and reverse  
strands

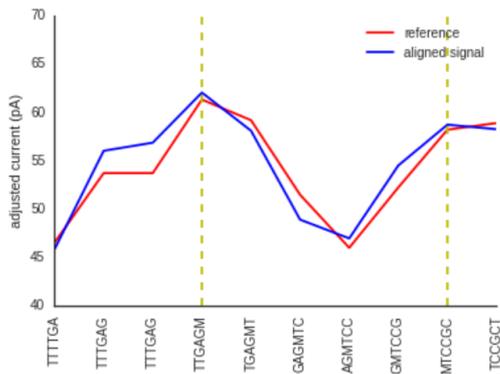
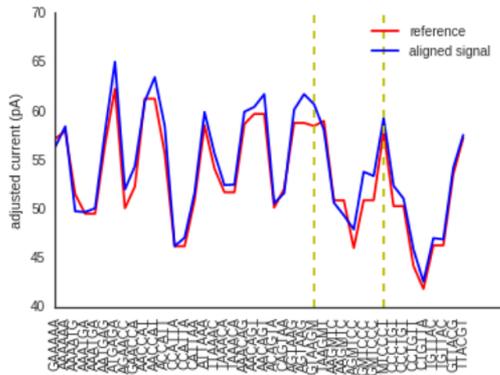
GATC

CTGCAG

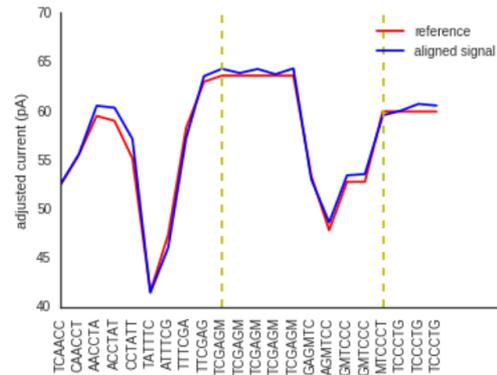
ACCACC

CCACN<sub>8</sub>TGAY/R  
TCAN<sub>8</sub>GTGG

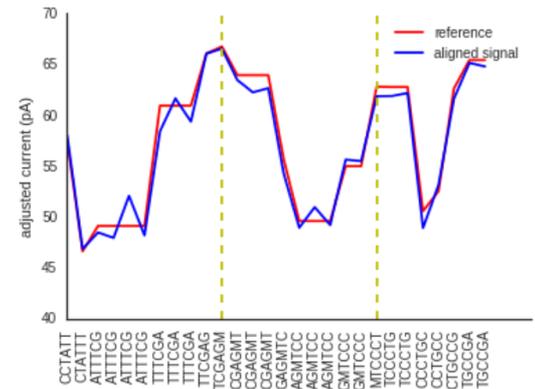
# Calling current (pA) differences, similar to PacBio



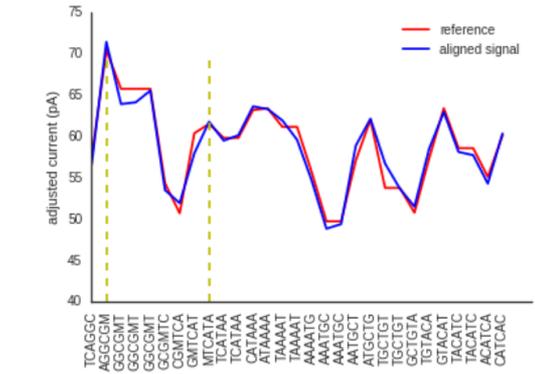
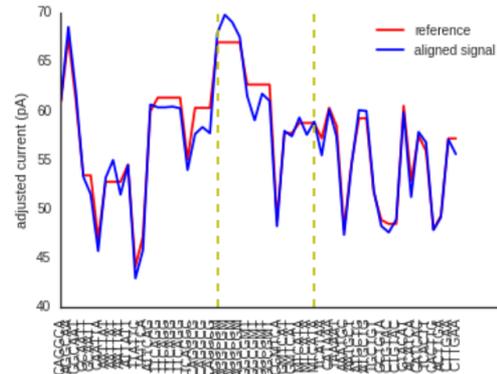
Reads aligned to same positions



||

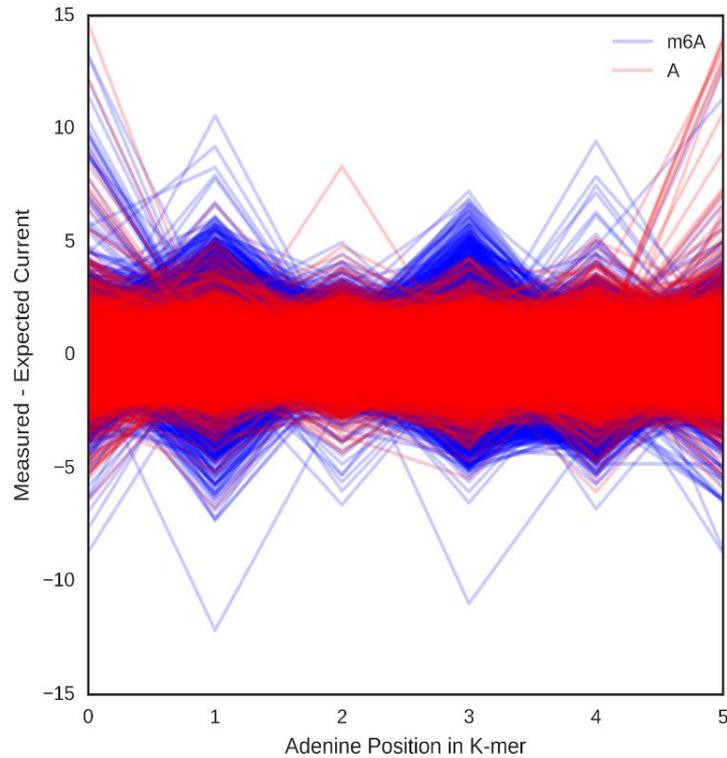


||

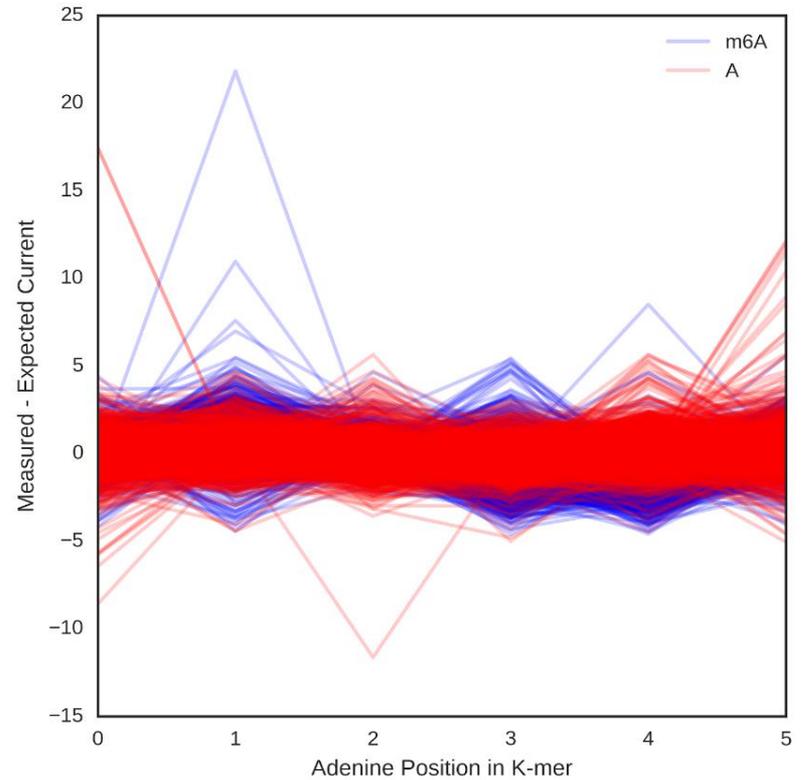


# Certain positions of the pore and more informative than others

Training run



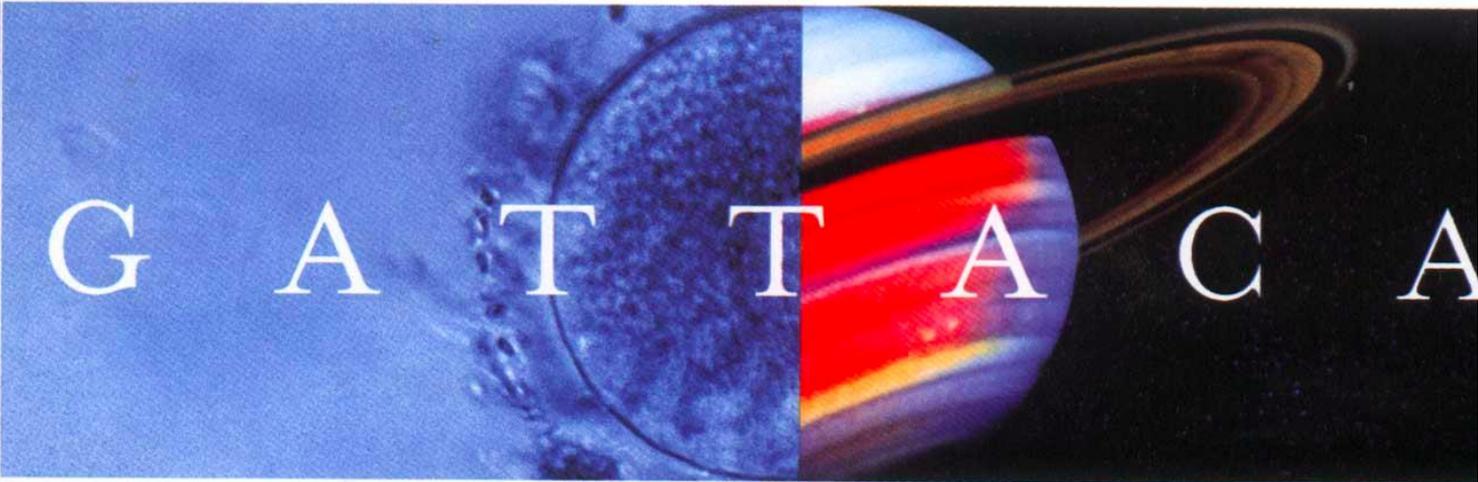
Test run





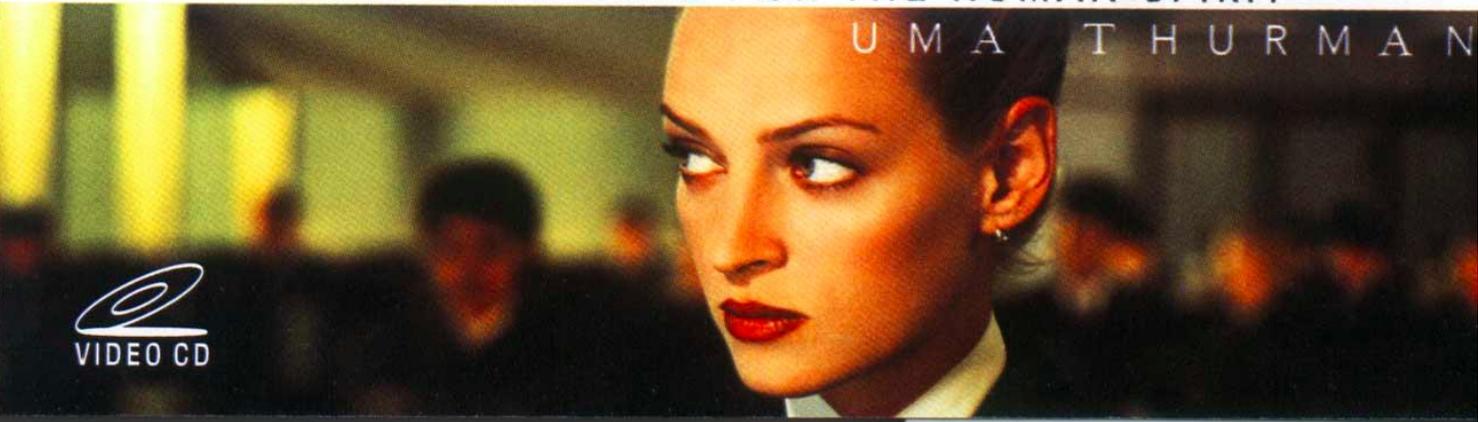


ETHAN HAWKE



GATTACA

THERE IS NO GENE FOR THE HUMAN SPIRIT



UMA THURMAN



VIDEO CD

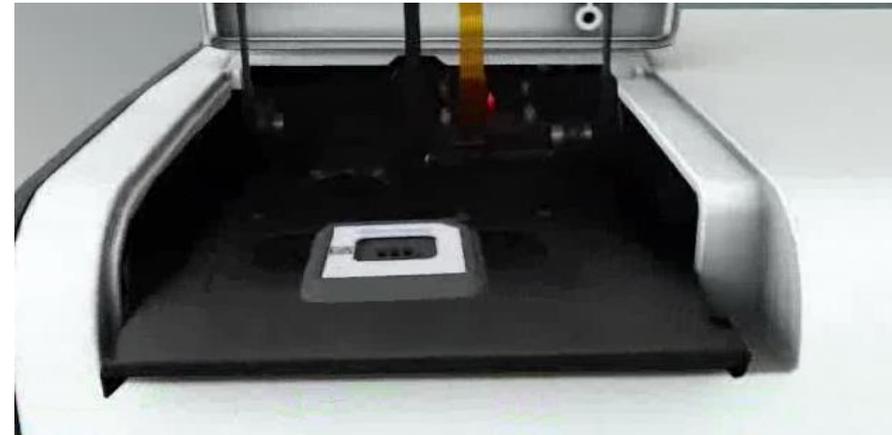
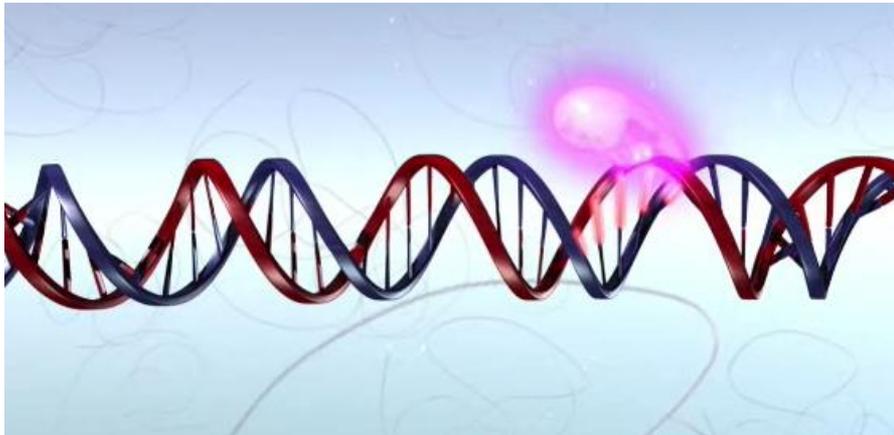
# Is a 2.6 minute genome possible?

## No today, but if the physics holds up...

Table 2: Nanopore and Nanochannel Sequencing Considerations

Parameter	DNA fragment (average bp)	Pore Speed (bp/s)	# nanopores	% of Pores Functional	transit time (seconds)	transit time (minutes)	run time (hours)	max # molecules / pore / run	% of time pores have DNA	actual # molecules/ pore/run	# of bases sequenced per device	Run Cost (\$)	\$ / Mb	\$ / Gb	Hours for 30X WGS of 3.1Gb	Model
Time	10,000	100	512	0.5	100	1.67	6	216	80%	172.8	442,368,000	\$ 1,000	\$ 2.26	\$ 2,260.56	1261.4	T1
	10,000	100	512	0.5	100	1.67	24	864	80%	691.2	1,769,472,000	\$ 1,000	\$ 0.57	\$ 565.14	1261.4	T2
	10,000	100	512	0.5	100	1.67	48	1728	80%	1382.4	3,538,944,000	\$ 1,000	\$ 0.28	\$ 282.57	1261.4	T3
Size	10,000	100	512	0.5	100	1.67	6	216	80%	172.8	442,368,000	\$ 1,000	\$ 2.26	\$ 2,260.56	1261.4	S1
	100,000	100	512	0.5	1000	16.67	6	21.6	80%	17.28	442,368,000	\$ 1,000	\$ 2.26	\$ 2,260.56	1261.4	S2
	1,000,000	100	512	0.5	10000	166.67	6	2.16	80%	1.728	442,368,000	\$ 1,000	\$ 2.26	\$ 2,260.56	1261.4	S3
Size & Time	10,000	100	512	0.5	100	1.67	6	216	80%	172.8	442,368,000	\$ 1,000	\$ 2.26	\$ 2,260.56	1261.4	S&T1
	100,000	100	512	0.5	1000	16.67	24	86.4	80%	69.12	1,769,472,000	\$ 1,000	\$ 0.57	\$ 565.14	1261.4	S&T2
	1,000,000	100	512	0.5	10000	166.67	48	17.28	80%	13.824	3,538,944,000	\$ 1,000	\$ 0.28	\$ 282.57	1261.4	S&T3
Pores	10,000	100	50000	0.5	100	1.67	6	216	80%	172.8	43,200,000,000	\$ 1,000	\$ 0.023	\$ 23.15	12.9	P&T1
	10,000	100	100000	0.5	100	1.67	6	216	80%	172.8	86,400,000,000	\$ 1,000	\$ 0.012	\$ 11.57	6.5	P&T2
	10,000	100	150000	0.5	100	1.67	6	216	80%	172.8	129,600,000,000	\$ 1,000	\$ 0.008	\$ 7.72	4.3	P&T3
Pores & Time	10,000	100	50000	0.5	100	1.67	6	216	80%	172.8	43,200,000,000	\$ 10,000	\$ 0.23	\$ 231.48	12.9	P&T1
	10,000	100	100000	0.5	100	1.67	24	864	80%	691.2	345,600,000,000	\$ 20,000	\$ 0.06	\$ 57.87	6.5	P&T2
	10,000	100	150000	0.5	100	1.67	48	1728	80%	1382.4	1,036,800,000,000	\$ 30,000	\$ 0.03	\$ 28.94	4.3	P&T3
Speed & Time	10,000	100	50000	0.5	100	1.67	6	216	80%	172.8	43,200,000,000	\$ 10,000	\$ 0.23	\$ 231.48	12.9	PS&T1
	10,000	1000	100000	0.5	10	0.17	24	8640	80%	6912	3,456,000,000,000	\$ 20,000	\$ 0.01	\$ 5.79	0.6	PS&T2
	10,000	10000	150000	0.5	1	0.02	48	172800	80%	138240	103,680,000,000,000	\$ 30,000	\$ 0.00	\$ 0.29	0.04	PS&T3

# Bionanogenomics - Irys System

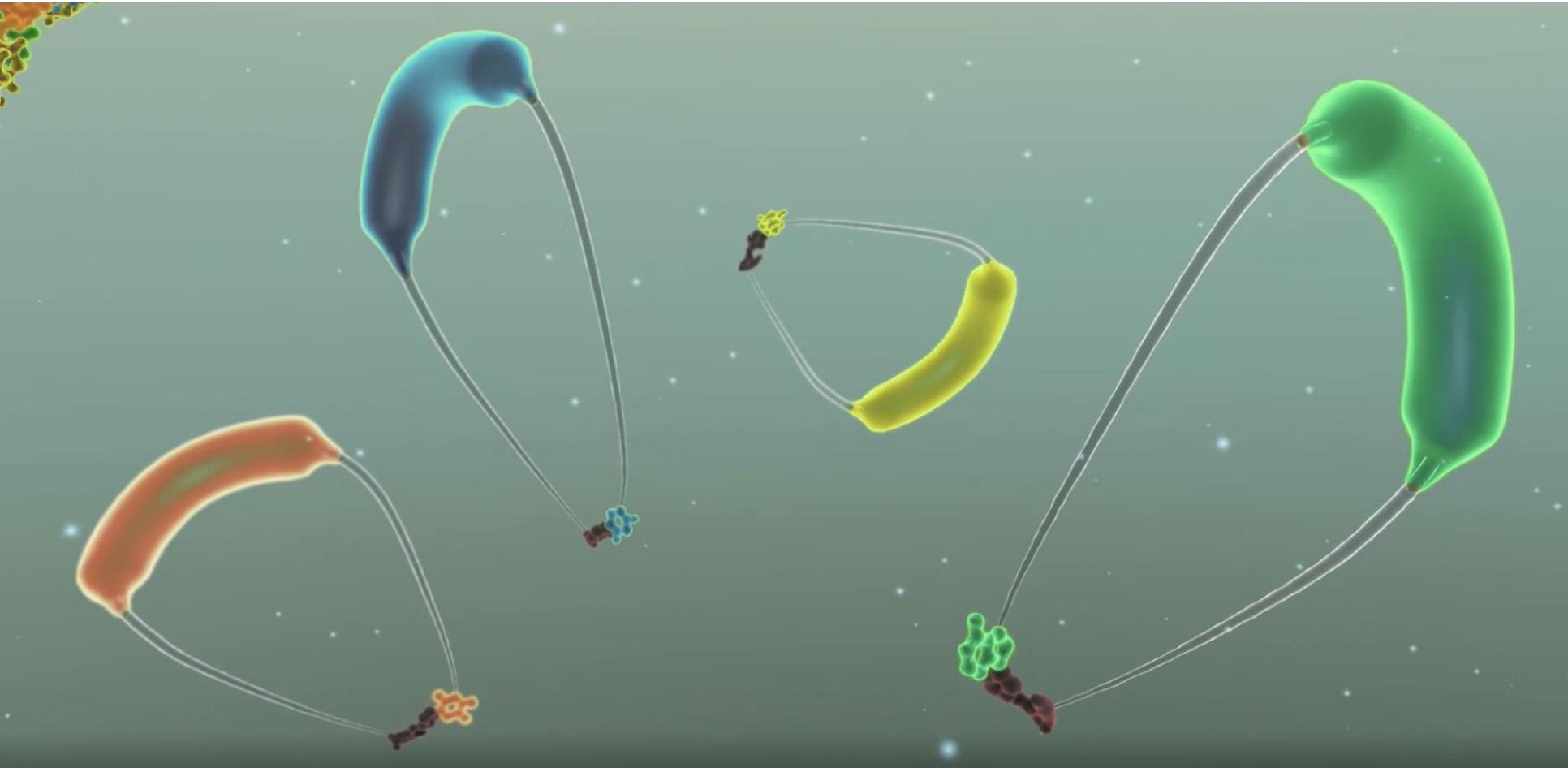


# The new Illumina iSeq100 can sequence in <6h.



# Emerging Technologies

# Roche's nanopore tech



Sequencing by eXpansion (SBX)

# The race for long is on

Longer and longer: DNA sequence of more than two million bases now achieved with nanopore sequencing.

Fri 4th May 2018

## Congratulations!

The first >2 Mb DNA read, achieved with nanopore sequencing

Matt Loose, Alex Payne, Nadine Holmes, Vardhman Rakyan & team, University of Nottingham, UK

May 2018

Long read  
club



Really very long reads  
indeed

<http://longreadclub.org/>

<https://nanoporetech.com/about-us/news/longer-and-longer-dna-sequence-more-two-million-bases-now-achieved-nanopore>

# News

10/31/2018

## BGI Unveils New High-Throughput Sequencing System.

Last week at the 13<sup>th</sup> International Conference on Genomics (ICG-13) in Shenzhen, China, BGI announced a new sequencing system based on its DNBseq™ Technology.

The newly unveiled **MGISEQ-T7** is the most powerful sequencing system from BGI's MGI subsidiary, with a daily output capability of 6Tb of data.

The MGISEQ-T7 is able to complete 60 human genomes in a single day, with essentially error-free sequencing from BGI's DNBseq sequencing technology.

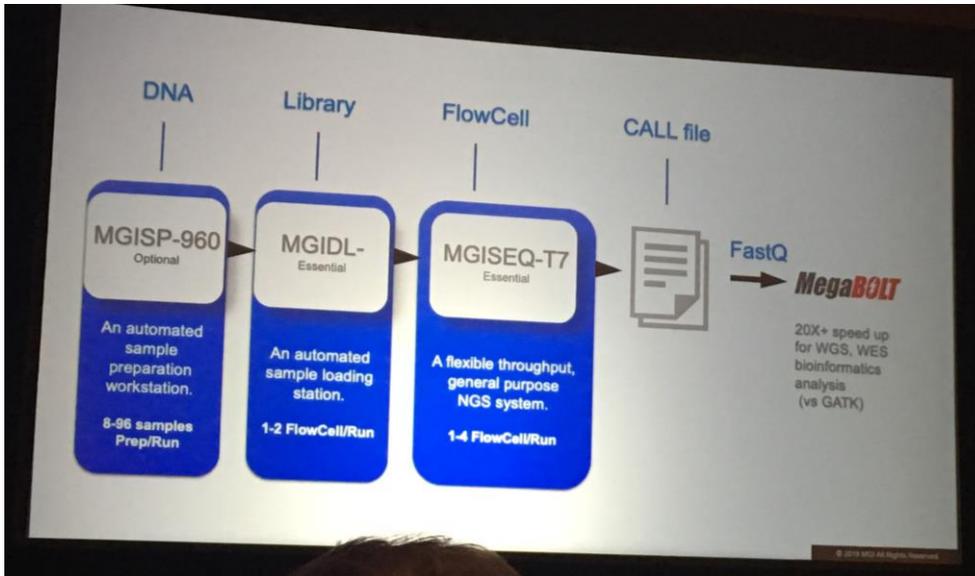




T-1000?



# BGI – NGS streets

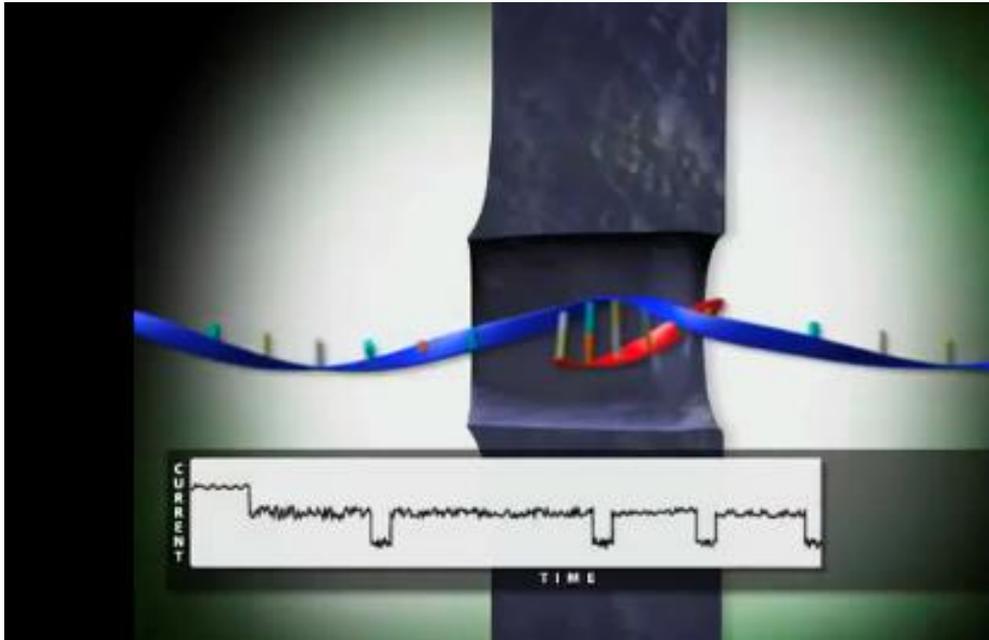
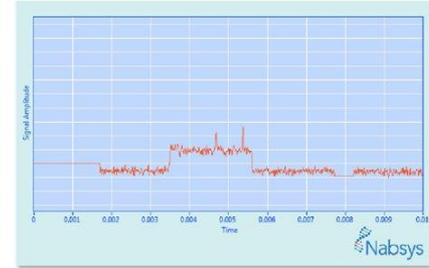
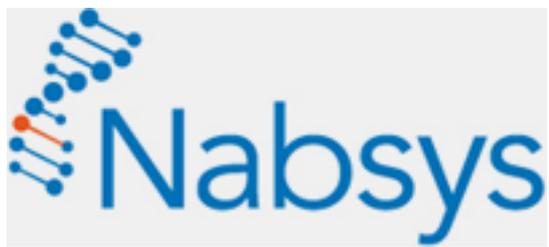


omics for All"

**1000 Sequencing labs x 1000 Genomes/day**

**100M Birth Rate + 7B Population Sequenced at 50yrs → 240M/yr**

© 2015 MGI All Rights Reserved.



## Hybridization -Assisted Nanopore Sequencing (HANS):

- 1 million bases per second
- Variable probe length can be used for HANS
- Long Reads (100kb)
- Single molecule



Single-atom labeling and then visualization with EM

- Long Reads (20kb)
- Single molecule

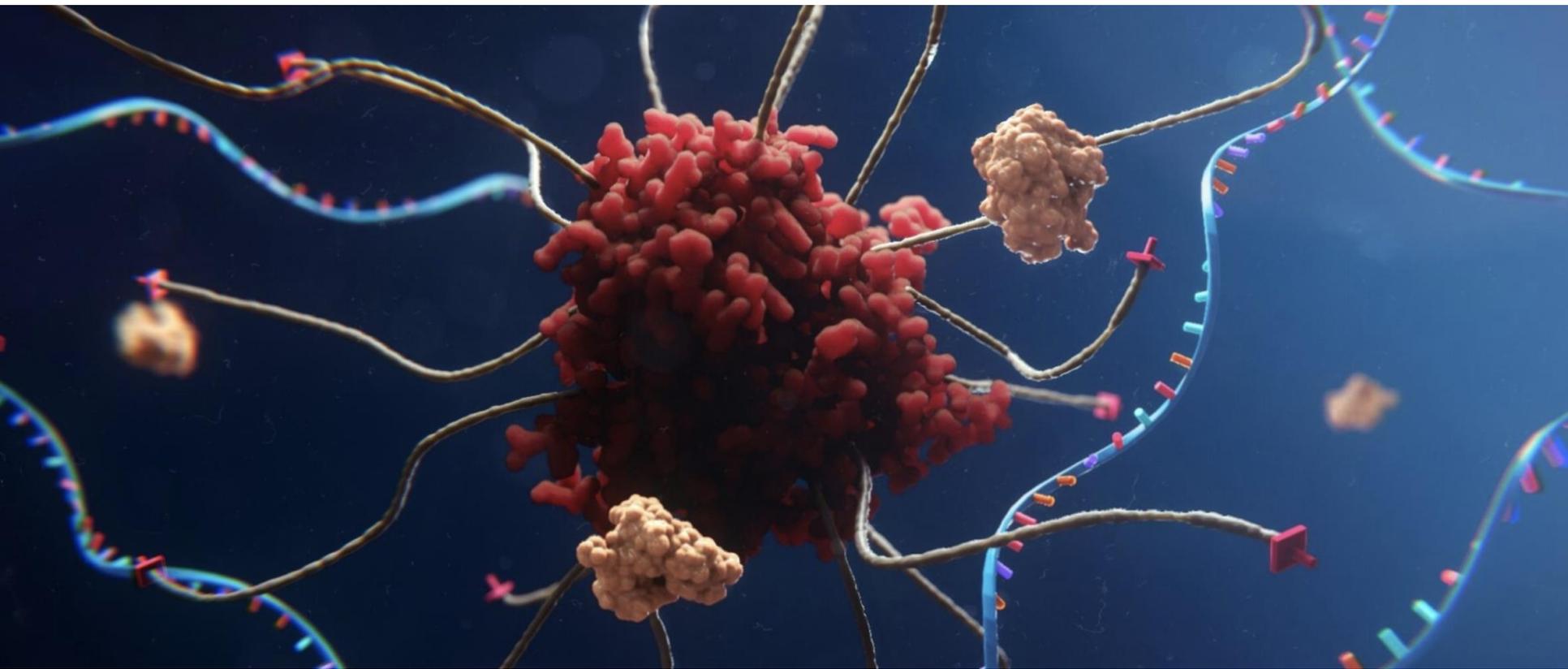
# GenapSys



(1M, 16M or 144M)



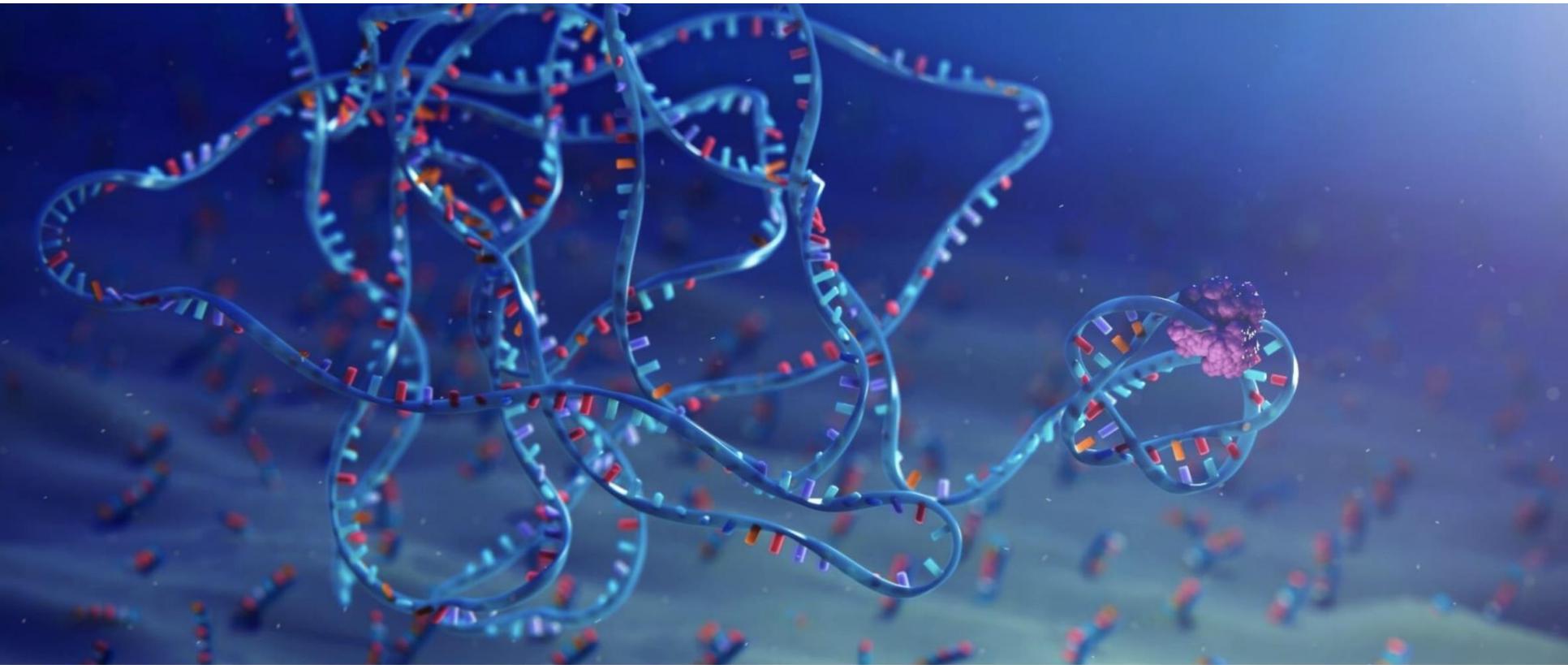
# New Kid on the Block: Element Biosciences



Element  
Biosciences

<https://www.elementbiosciences.com/technology>

# Surface-anchored amplification sites



Element  
Biosciences

<https://www.elementbiosciences.com/technology>

# Each Platform has various sources of noise, and thus Error

- De-Phasing
  - Lagging strand dephasing from incomplete extension
  - Leading strand dephasing from over-extension
- Dark Nucleotides
- Polymerase errors ( $10^{-5}$  to  $10^{-7}$ )
- Single molecule challenges
  - High noise
  - Polymerase “wiggling” from tail
- Platform-specific errors
  - Illumina more likely to have error after ‘G’
  - PCR-based methods miss GC- and AT-rich regions



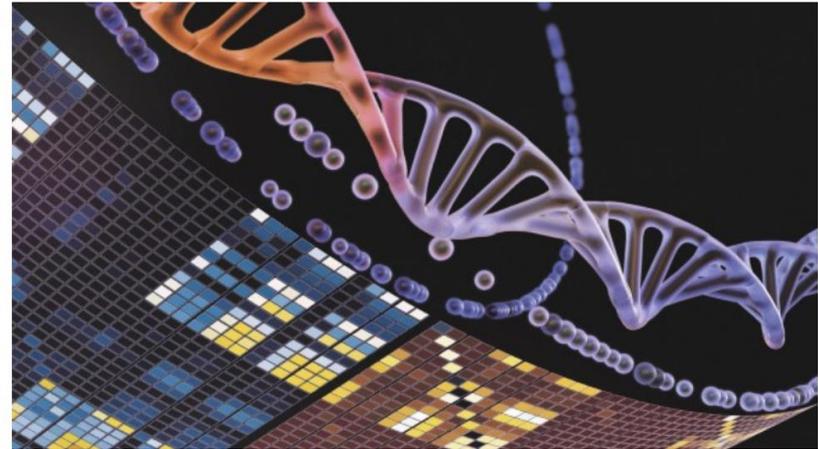
# How can we be sure we sequenced correctly?

nature > nature biotechnology > collection

COLLECTION | 09 SEPTEMBER 2021

## Sequencing Quality Control 2

This Web Collection presents the results of the Sequencing Quality Control 2 (SEQC2) project that sought to evaluate quality-control metrics and human, bacterial and metagenomic reference materials and datasets for next-generation sequencing (NGS) in both regulatory... [show more](#)



<https://www.nature.com/collections/seqc2>

# Test them!

nature biotechnology

Explore content ▾

About the journal ▾

Publish with us ▾

Subscribe

[nature](#) > [nature biotechnology](#) > [articles](#) > article

Article | [Published: 09 September 2021](#)

## Performance assessment of DNA sequencing platforms in the ABRF Next-Generation Sequencing Study

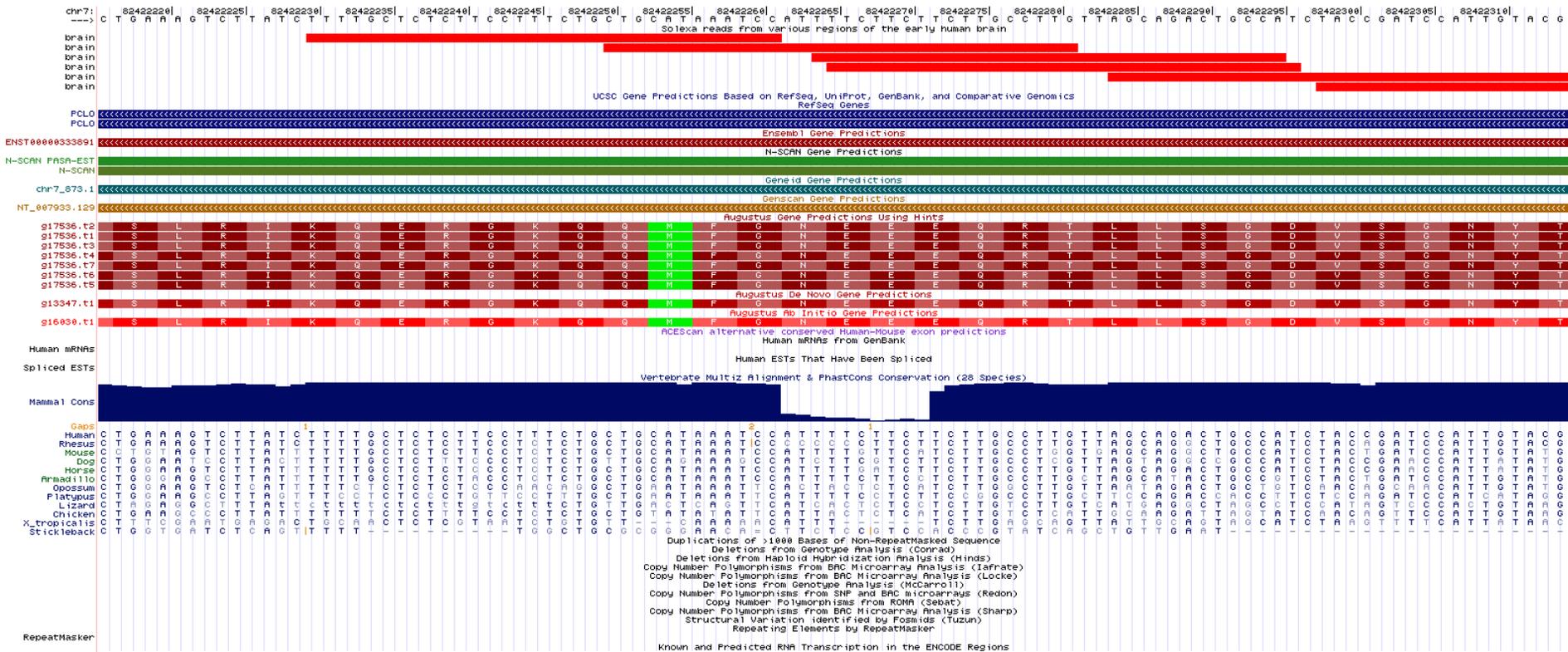
[Jonathan Foox](#), [Scott W. Tighe](#), ... [Christopher E. Mason](#) 

+ Show authors

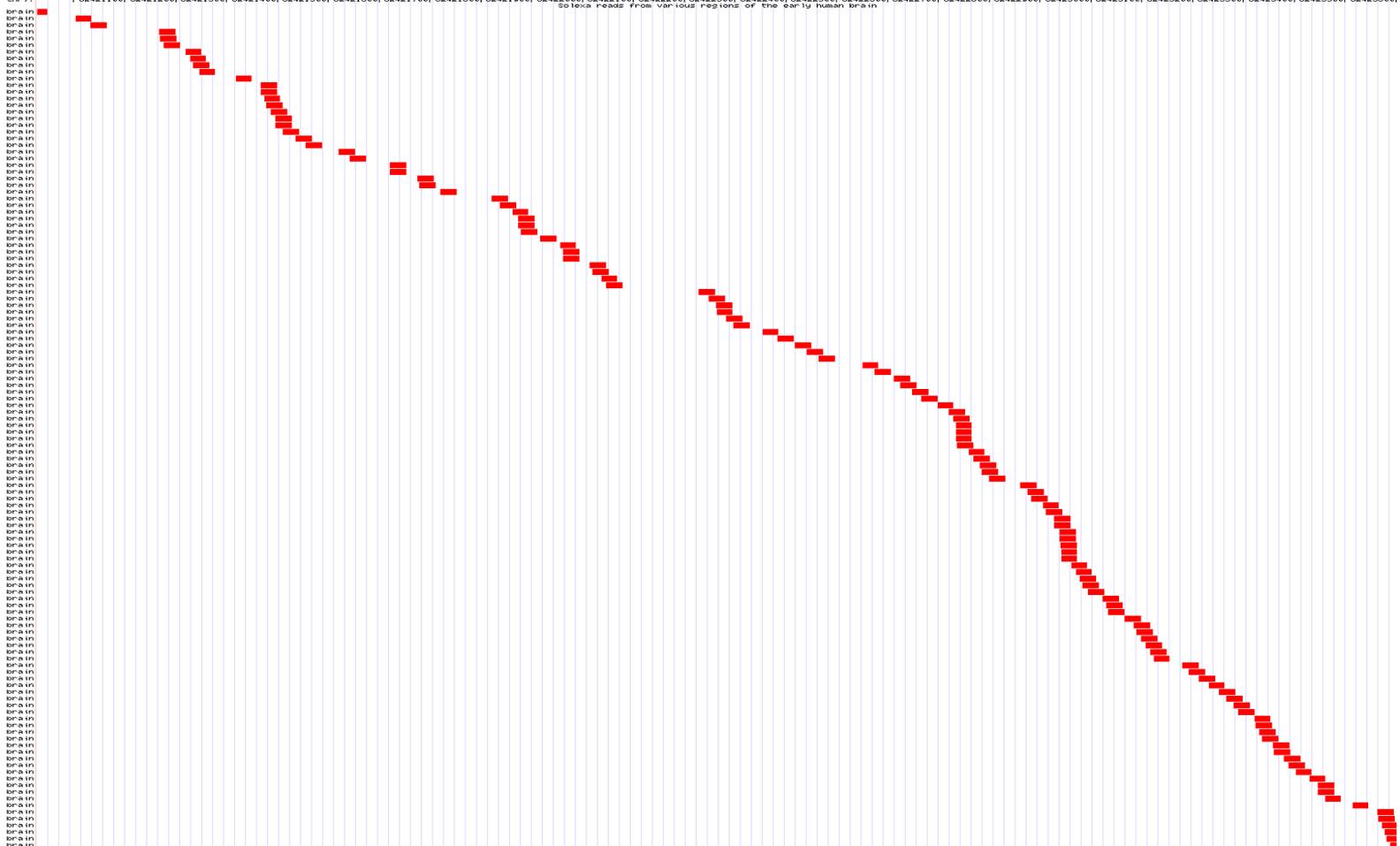
<https://www.nature.com/articles/s41587-021-01049-5>

What do you do with the reads?

# Alignment to the genome



chr7:1 | 82421100 82421200 82421300 82421400 82421500 82421600 82421700 82421800 82421900 82422000 82422100 82422200 82422300 82422400 82422500 82422600 82422700 82422800 82422900 82423000 82423100 82423200 82423300 82423400 82423500 82423600  
50bp bins from various regions of the chr7 human brain



UCSC Gene Predictions Based on RefSeq, UniProt, GenBank, and Comparative Genomics

RefSeq Gene

Ensembl Gene Predictions

N-SCRN Predictions

Gene3D Gene Predictions

GenScan Gene Predictions

Augustus Gene Predictions Using Hints

Augustus De Novo Gene Predictions

Augustus Ab Initio Gene Predictions

NCSCan alternative conserved RepeatMasker exon predictions

Human mRNAs

Human ESTs That Have Been Spliced

Spliced ESTs

Vertebrate Multiz Alignment & PhyloP Conservation (28 Species)

Mammal Cons

Rhesus

Mouse

Dog

Hop 50

Prad 110

Genus 10

Fatopus

120

Chickon

Xtropic 100

St ck Teback

RepeatsMasker

Known and Predicted RNA Transcription in the ENCODE Regions

Structural Variations from the 1000 Genomes Project

Copy Number Polymorphisms from SNP and BAC Microarrays (Redon)

Copy Number Polymorphisms from BAC Microarray Analysis (Locke)

Deletions from Rapid Hybridization Analysis (Hinds)

Deletions from Genotype Analysis (Cibulka)

Copy Number Polymorphisms from SNP and BAC Microarrays (Redon)

Copy Number Polymorphisms from SNP and BAC Microarrays (Locke)

Structural Variations from the 1000 Genomes Project

RepeatsMasker

Known and Predicted RNA Transcription in the ENCODE Regions

# The reads: FASTQ

The most common format is FASTQ, based off the FASTA data format:

```
>SequencedID
```

```
CGTAGTCTATATATGCGCGAATGCGTA
```

**But....**

FASTQ also includes quality information:

```
@Sample_Info
```

```
CCTTGCTGCC
```

```
+
```

```
3.6;#$!>><
```

# Understanding FASTQ

For Illumina, sequences have an ID:

@HWUSI-EAS100R:6:73:941:1973#0/1

HWUSI-EAS100R	the unique instrument name
6	flowcell lane
73	tile number within the flowcell lane
941	'x'-coordinate of the cluster within the tile
1973	'y'-coordinate of the cluster within the tile
#0	index number for a multiplexed sample (0 for no indexing)
/1	the member of a pair, /1 or /2 (paired-end or mate-pair reads only)

# Understanding Quality Scores

Q-values are the probability (p) of a base being incorrect. From Sanger sequencing:

$$Q_{\text{value}} = -10 \log_{10} p$$

---

So, if your  $p=0.1$ , then  $Q_{\text{value}} = (-10 \log_{10}(0.1))$   
 $= (-10(-1)) = 10$

---

If your  $p=0.01$ , then  $Q_{\text{value}} = (-10 \log_{10}(0.01))$   
 $= (-10(-2)) = 20$

---

If  $p=0.001$ , then  $Q_{\text{value}} = (-10 \log_{10}(0.001))$   
 $= (-10(-3)) = 30$



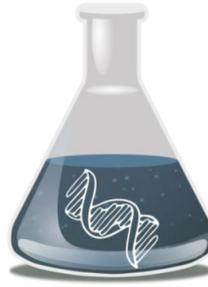




<https://jimb.stanford.edu/giab/>

# GENOME IN A BOTTLE

*Our mission is to provide the authoritative characterization of human genomes.*



[OVERVIEW](#)

[WORKSHOPS](#)

[NEWS](#)

[RESOURCES](#)

[GOOGLE GROUP](#)

[ANALYSIS TEAM](#)

## Reference Materials and Data

The Genome in a Bottle Consortium has selected several genomes to produce and characterize as reference materials. The National Institute of Standards and Technology (NIST) is developing NIST Reference Materials from these genomes, which are DNA extracted from a large homogenized growth of B lymphoblastoid cell lines from the Coriell Institute for Medical Research. Note that there may be small differences between the NIST DNA and the Coriell DNA since they come from different growths of cells, though we do not expect these differences to be large for most applications.

The NIST Reference Materials available and planned are listed below, along with links to their data.

A description of data generated by GIAB for all the genomes below is published [here](#), and characterization of small variants is published [here](#). Ongoing work to characterize more difficult variants and regions is announced in the [GIAB Analysis Team google group](#).

# Metagenome in a bottle



PRODUCTS ▾

SERVICES ▾

HOW TO ORDER ▾

RESOURCES ▾

ABOUT ▾



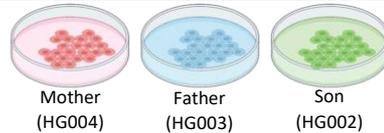
## COMPARISON TABLE

Catalog #	Product	Size
D6300	ZymoBIOMICS Microbial Community Standard	10 Preps
D6305	ZymoBIOMICS Microbial Community DNA Standard	200 ng
D6310	ZymoBIOMICS Microbial Community Standard II (Log Distribution)	10 Preps
D6311	ZymoBIOMICS Microbial Community DNA Standard II (Log Distribution)	220ng/20µl
D6320	ZymoBIOMICS Spike-in Control I (High Microbial Load)	25 Preps
D6321	ZymoBIOMICS Spike-in Control II (Low Microbial Load)	25 Preps
D6322	ZymoBIOMICS HMW DNA Standard	5000 ng
D6323	ZymoBIOMICS Fecal Reference with TruMatrix™ Technology	10 preps
D6331	ZymoBIOMICS Gut Microbiome Standard	10 preps

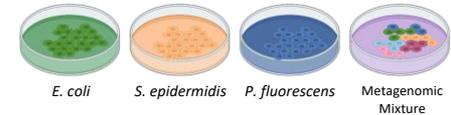
<https://www.zymoresearch.com/collections/zymbiomics-microbial-community-standards>

# ABRF-NGS Study Overview

Fully consented HapMap  
Ashkenazi Trio (NIST RM 8392)



Bacterial Genomes  
(ATCC MSA-3001)



## Whole Genome

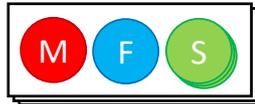
HiSeq2500



HiSeq4000



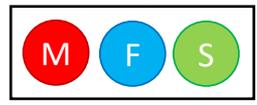
HiSeqX10



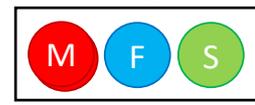
NovaSeq (2x150 and 2x250)



BGI-SEQ 500



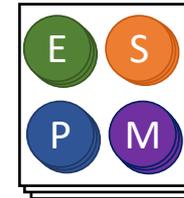
MGI-SEQ 2000



MiSeq



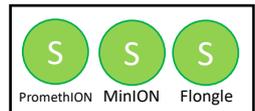
Ion PGM



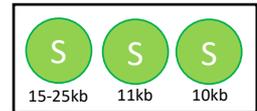
Ion S5



Oxford Nanopore



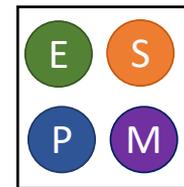
PacBio CCS



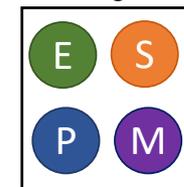
Genapsys



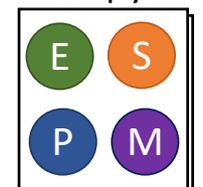
MinION



Flongle



Genapsys

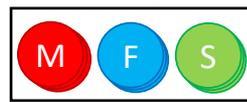


## Targeted Exome

Ion Proton

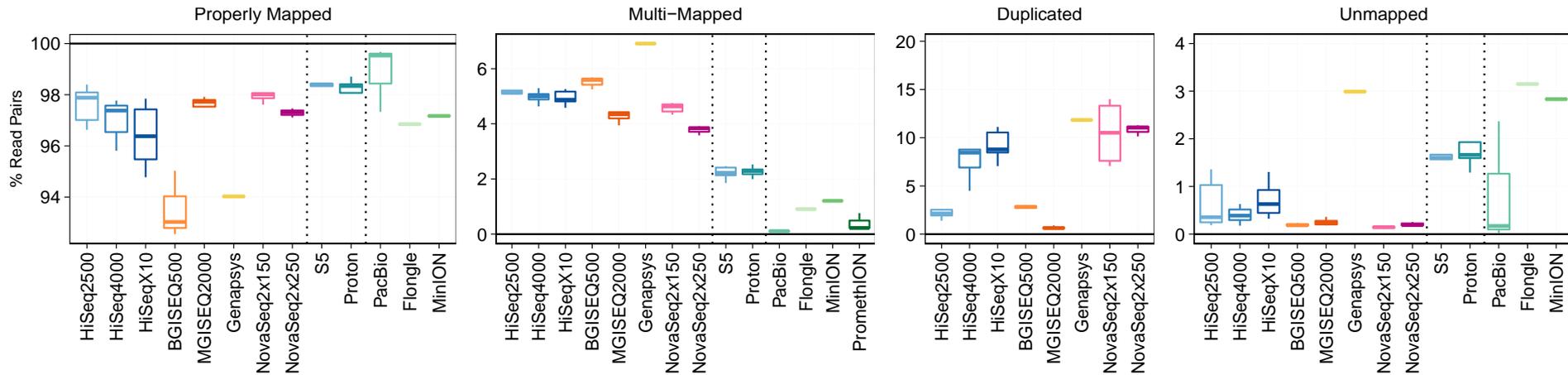
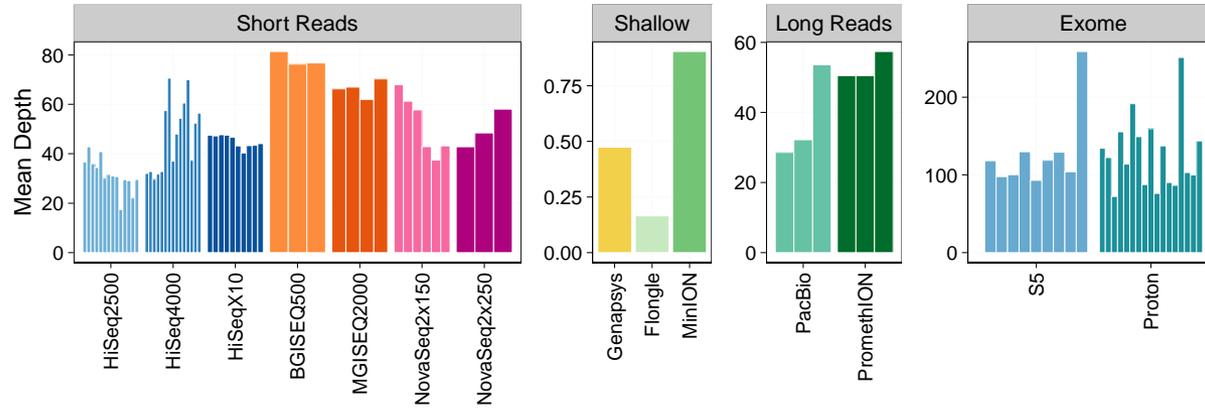


Ion S5



<https://www.nature.com/articles/s41587-021-01049-5>

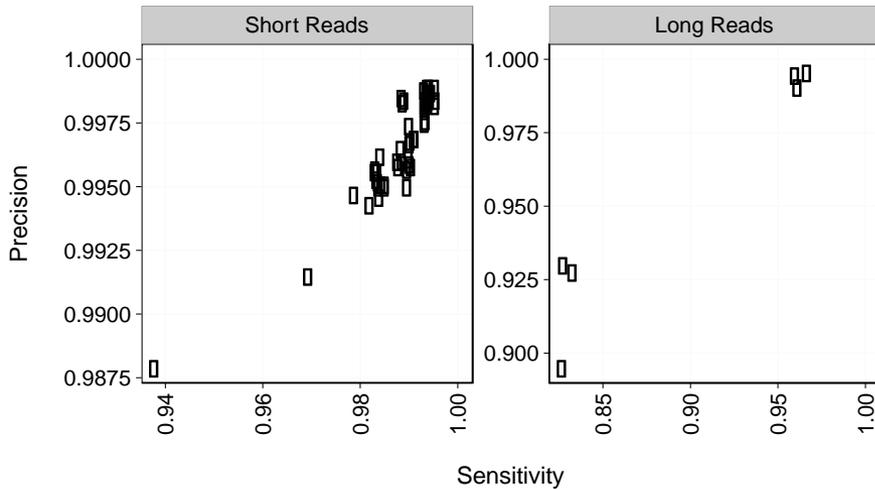
# Sequencing Depth and Mapping Efficiencies (Human Samples)



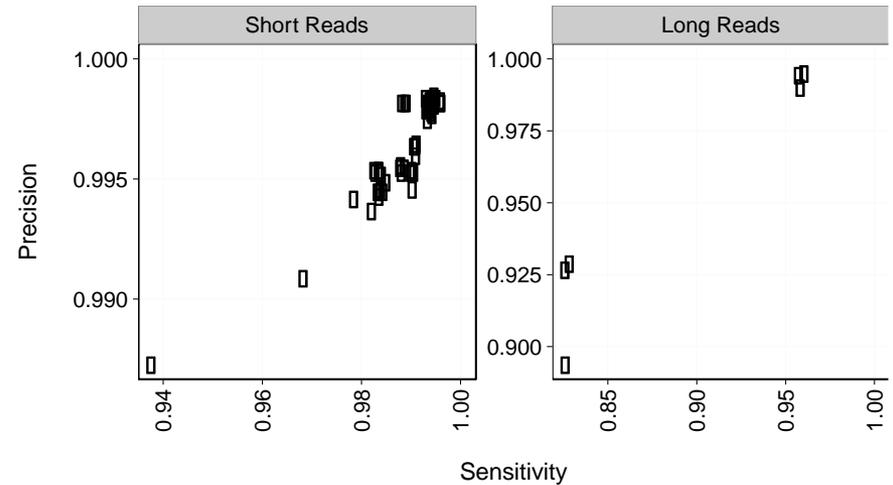


# Variant Calling in... Medically Relevant Genes

CLINVAR Variants



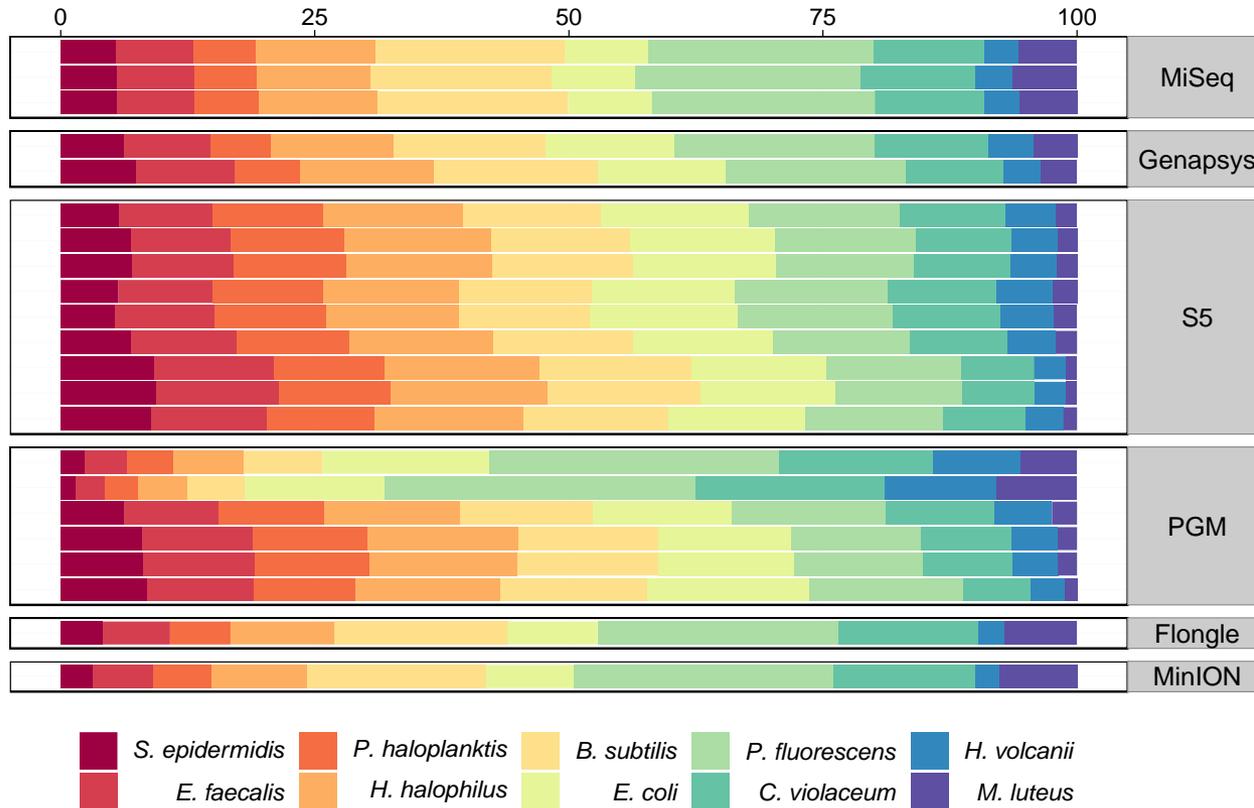
OMIM Variants



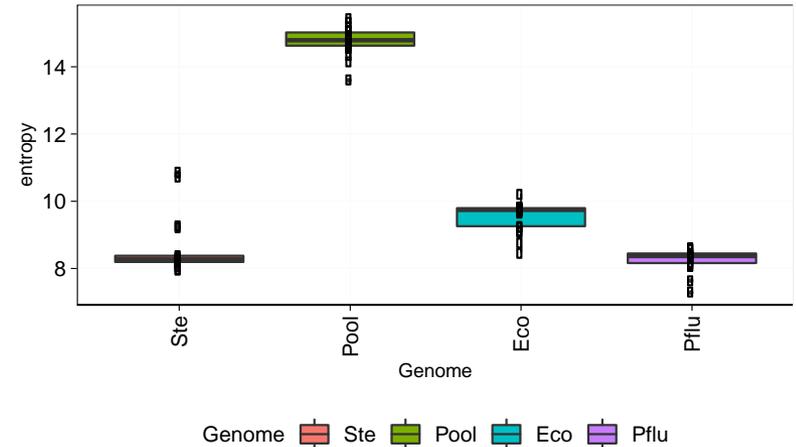
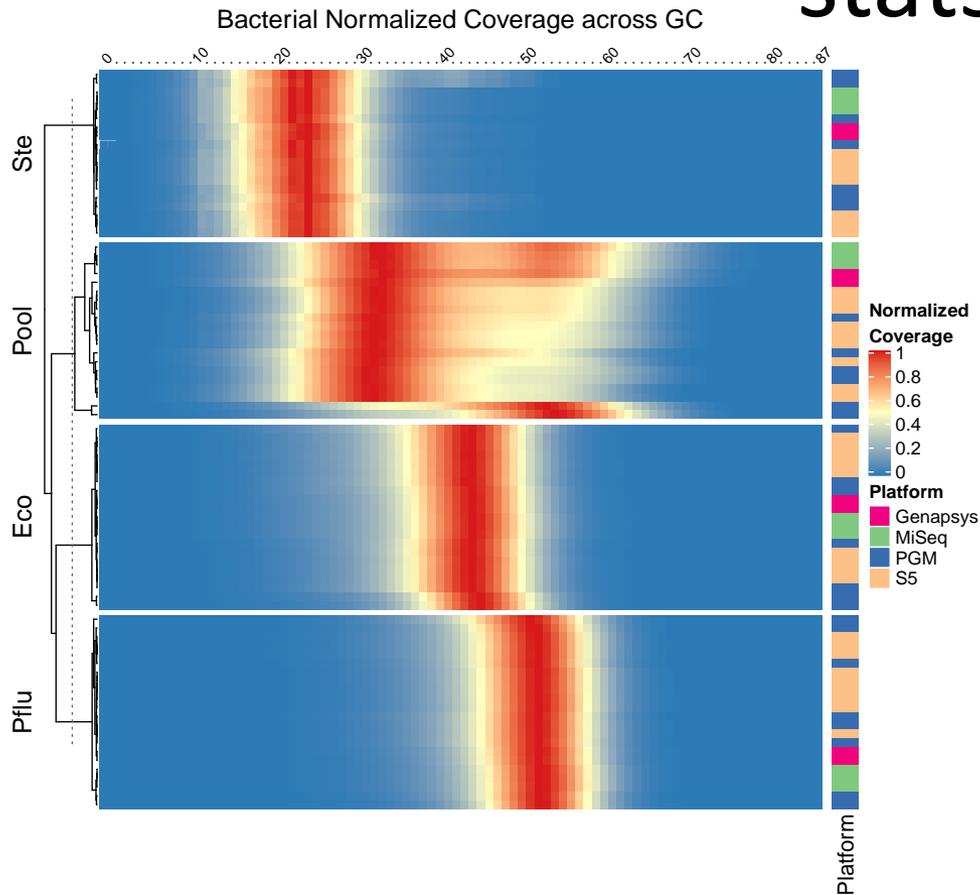
- ▣ HiSeq2500    ▣ NovaSeq2x150    ▣ MGISEQ2000
- ▣ HiSeq4000    ▣ NovaSeq2x250    ▣ PacBio
- ▣ HiSeqX10    ▣ BGISEQ500    ▣ Nanopore



# Metagenomic Samples are relatively similar



# Mixed samples mirror the pooled GC stats



# All Code and Data Available!

 <b>jfoox</b> Create variantAllele_GTtoMatrix.py	502bf0e 13 days ago	 45 commits
 Rmds	Mismatch Rates	13 days ago
 SLURM	Create strelka2.slurm	13 days ago
 bin	Create variantAllele_GTtoMatrix.py	13 days ago
 README.md	Update README.md	13 days ago
 longReadLinks.txt	Create longReadLinks.txt	13 days ago

README.md 

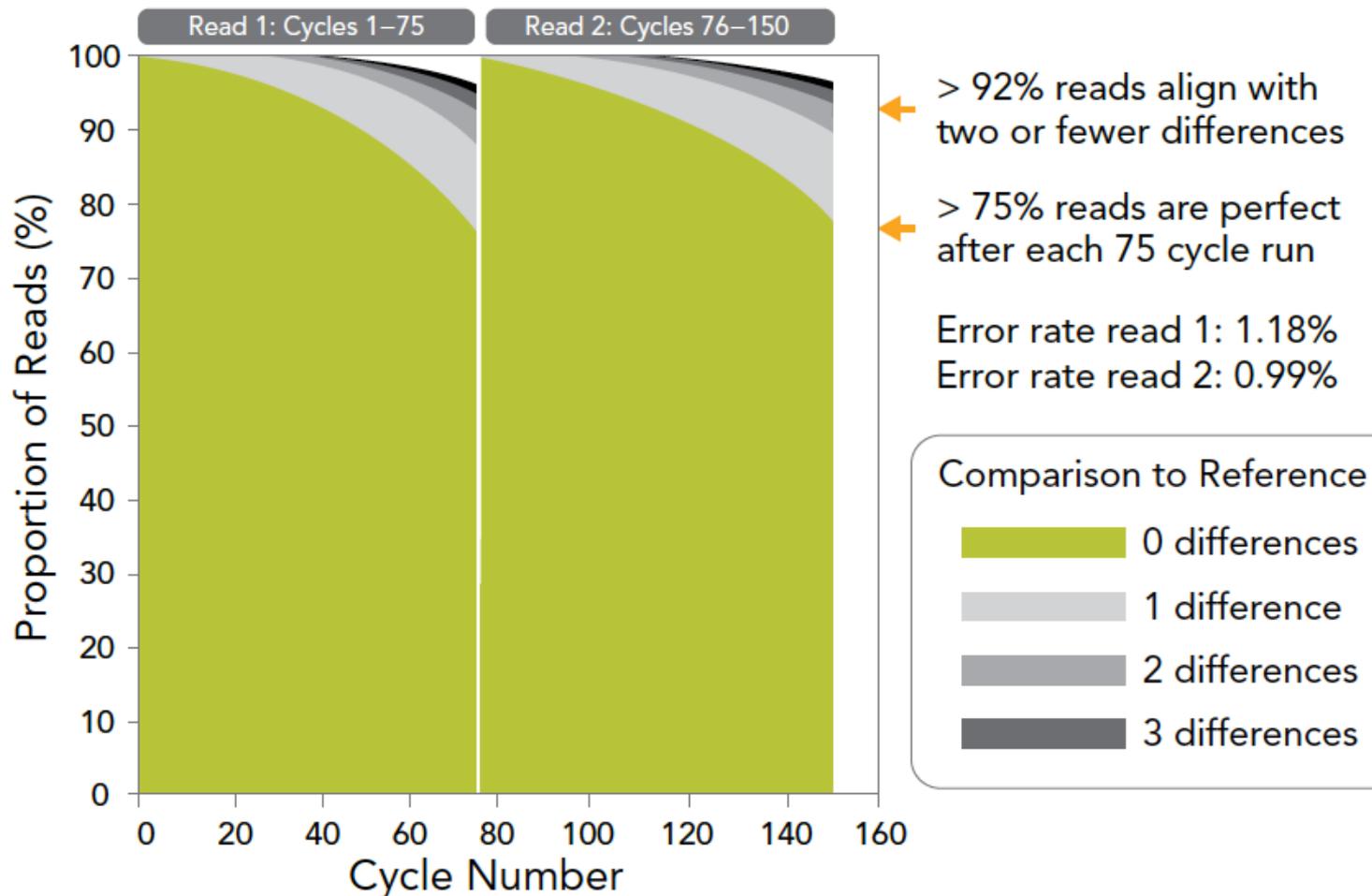
## ABRF NGS Phase II

Analysis and figure generation code for the ABRF NGS Phase II Study on DNA-seq reproducibility. This repository includes scripts to run heavy lifting such as alignment and variant calling (SLURM), shell scripts to do post-processing calculations (bin), and R scripts used to create figures (Rmds).

<https://github.com/jfoox/abrfngs2>

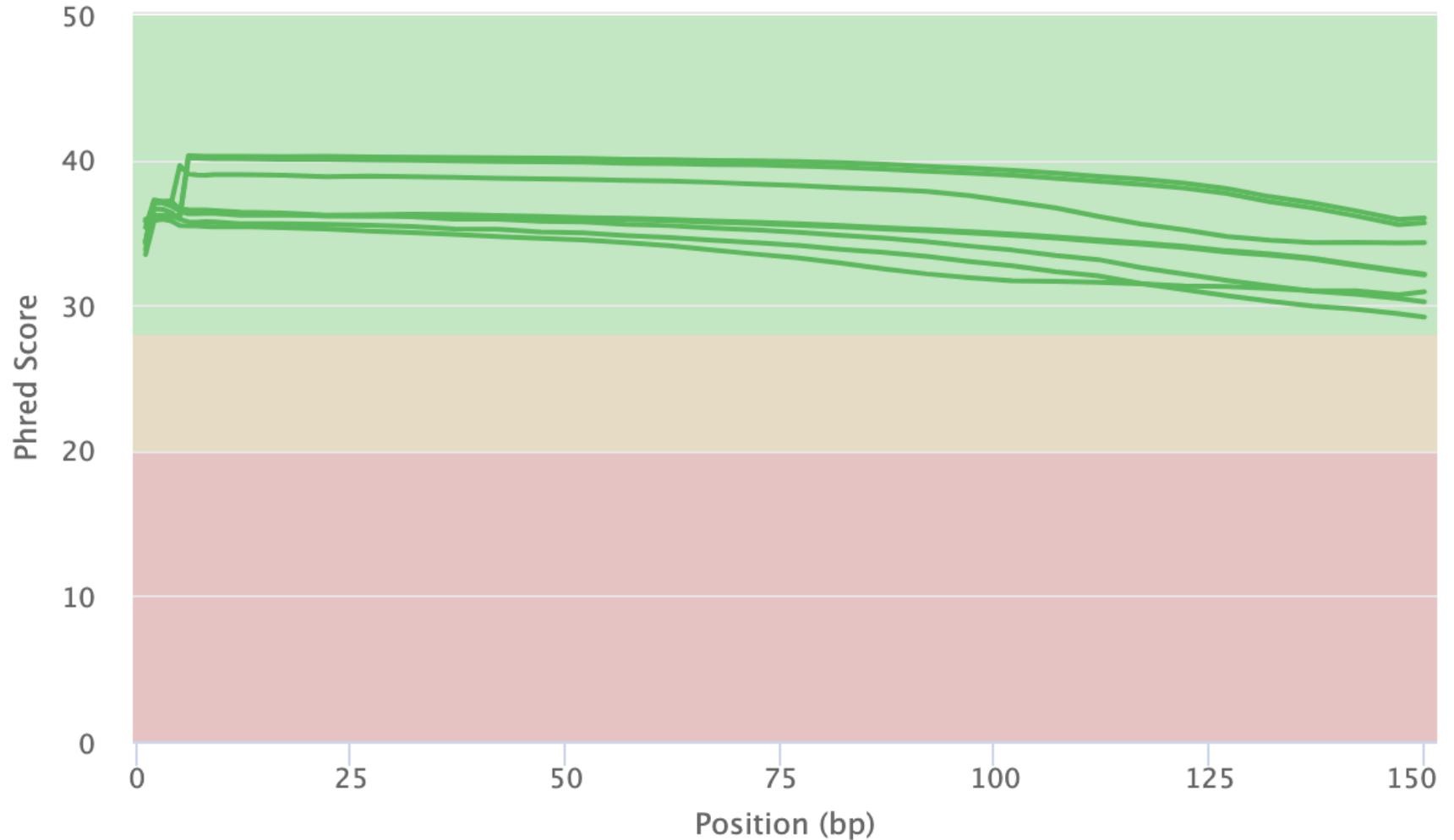
<https://www.nature.com/articles/s41587-021-01049-5>

# Many platforms are cycle-dependent on error rate - ILMN



# Element looks better so far

FastQC: Mean Quality Scores



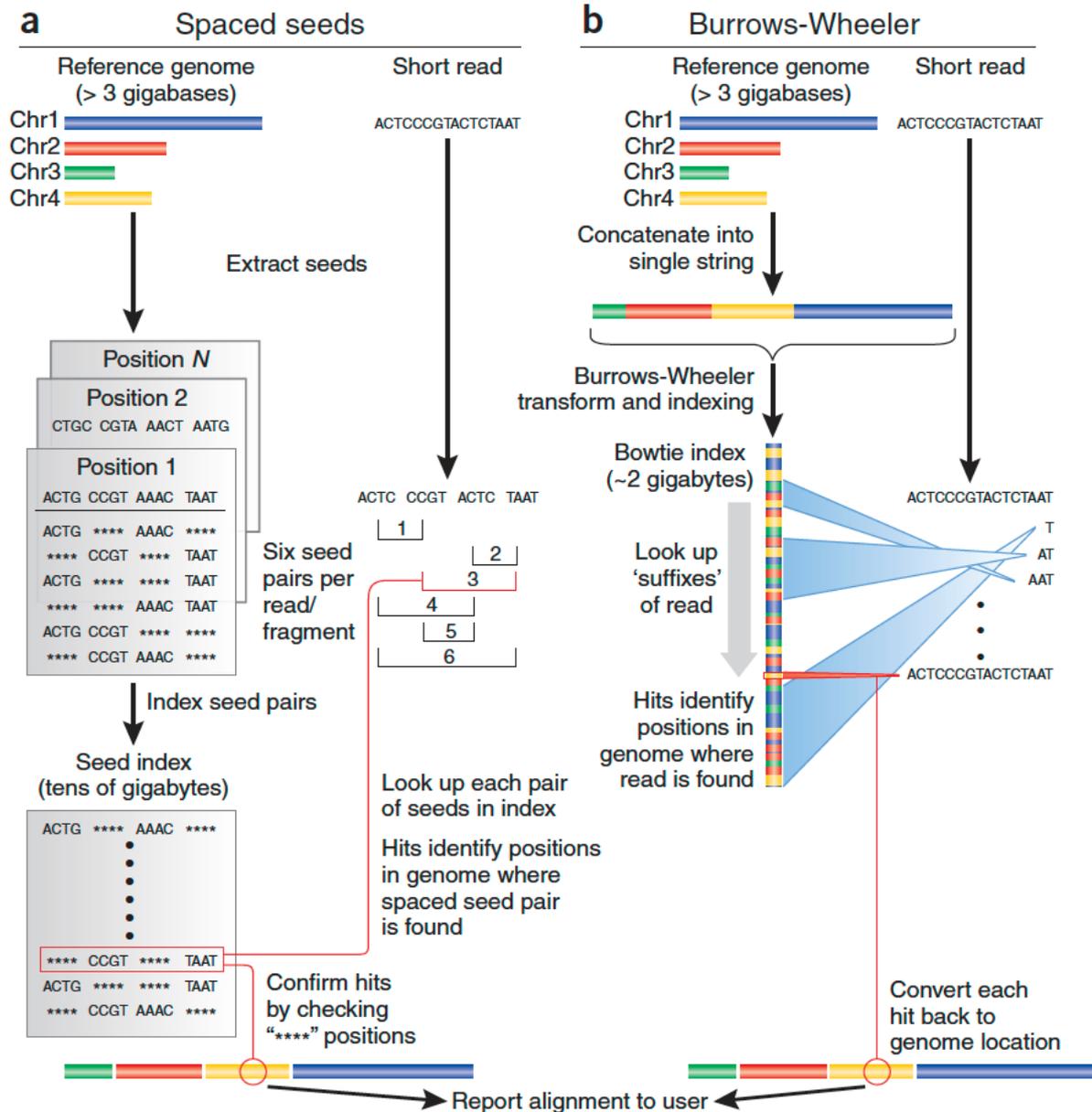
# Many Options for Alignment - 2009

	MAQ	ELAND	SOAP	BFAST	Bowtie	SHRiMP	Rmap	SeqMap	Novocraft
<b>Algorithm Parameters</b>									
Version	0.71	1.1	1.11	0.1.11	0.9.8	1.1.0	0.41	1.0.8	1.06
SNP-calls	✓	-	✓	-	-	✓	-	-	-
Uses Quality Scores	✓	-	-	✓	✓	✓	✓	-	✓
Indels	PE only	PE only	✓	✓	-	✓	-	✓	-
Splicing	-	-	-	-	-	-	-	-	-
Paired-End	✓	✓	✓	✓	-	-	-	-	✓
Threading	-	✓	✓	✓	✓	-	-	-	✓
Max # Mismatches (*in Seed)	3*	2*	5	-	3*, or UD	-	-	2	7
Default Seed Size	10	32	-	-	28	-	-	-	-
Max Input Length	63	-	60	-	-	-	64	-	-
5' Read Trimming	-	✓	-	-	✓	-	-	-	-
3' Read Trimming	✓	✓	✓	-	✓	-	-	-	✓
Methylation Alignment	-	-	-	✓	-	-	-	-	-
Repeats/Adaptor Removal	✓	✓	-	✓	✓	-	-	-	✓
Strand-specific search	-	-	✓	-	-	-	-	✓	-
<b>Platforms</b>									
ABI SOLiD	✓		✓	✓	✓	✓			
Illumina GA	✓	✓	✓	✓	✓	✓	✓	✓	✓
Roche 454					✓	✓			
Helicos Heliscope		✓	✓					✓	

# Many Options for Alignment - 2022

- Bfast
- BioScope
- Bowtie
- BWA
- CLC bio
- CloudBurst
- Eland/Eland2
- GenomeMapper
- GnuMap
- Karma
- MAQ
- MOM
- Mosaik
- MrFAST/MrsFAST
- NovoAlign
- PASS
- PerM
- RazerS
- RMAP
- SSAHA2
- Segemehl
- SeqMap
- SHRiMP
- Slider/SliderII
- SOAP/SOAP2
- Srprism
- Stampy
- vmatch
- ZOOM
- .....

# Many common methods are BW-based





# Burrows-Wheeler Transformation (BWT)

- First discovered in 1983 by Wheeler at AT&T Bell Labs
- Used for compression in 1994.
- First implemented for aligners with “Bowtie”  
Ben Langmead, Cole Trapnell, Mihai Pop,  
and Steven Salzberg
- Allows for fast searching with a small memory footprint

<http://bio-bwa.sourceforge.net/>

Li H. and Durbin R. “Fast and accurate short read alignment with Burrows-Wheeler transform.” (2009)  
*Bioinformatics*, 25, 1754-60.

Burrows M, Wheeler DJ. “A Block Sorting Lossless Data Compression Algorithm.” Technical Report 124. Palo Alto, CA: Digital Equipment Corporation; 1994.

Questions?