

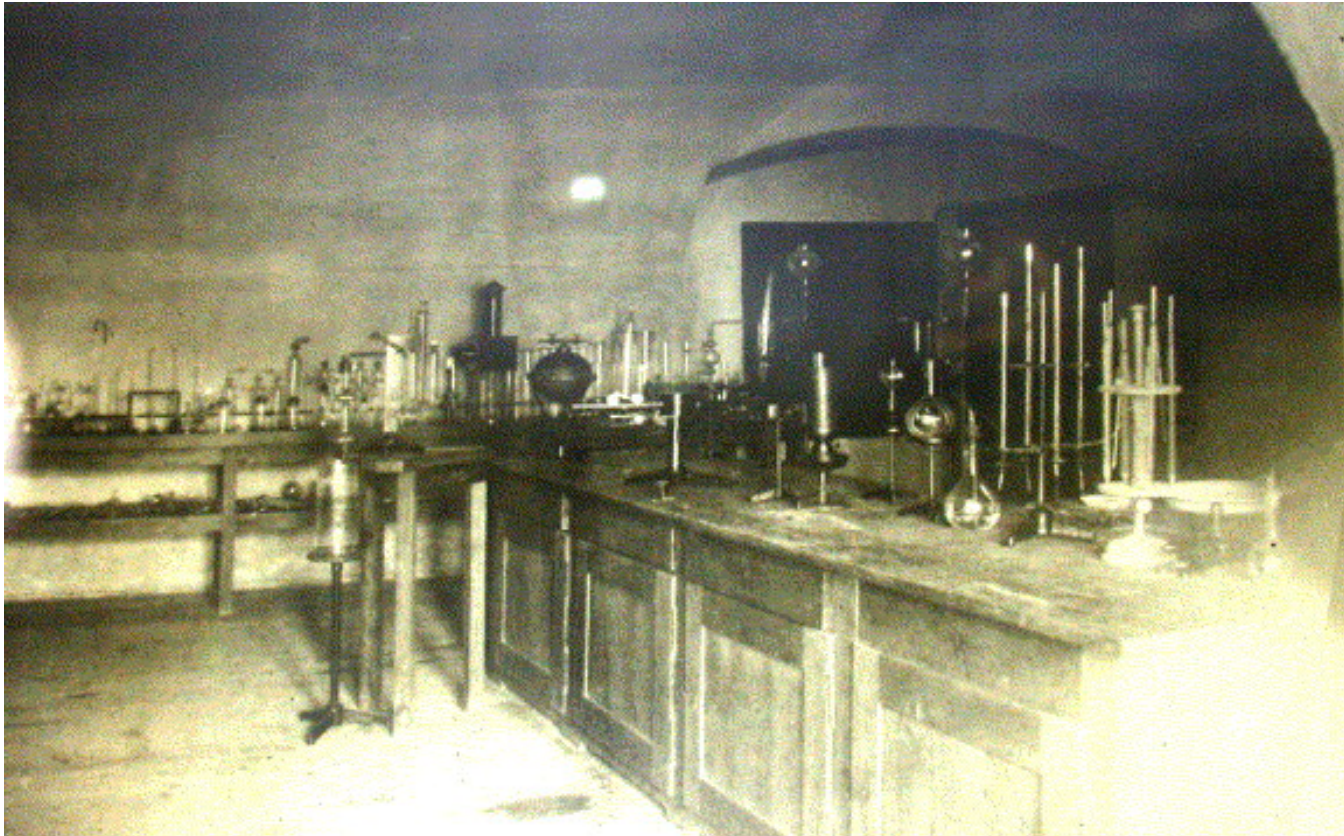
# Principles of sequencing: DNA, RNA

GBIO0002

2018-11-13

Tina O'Grady

# 1869: discovery of DNA



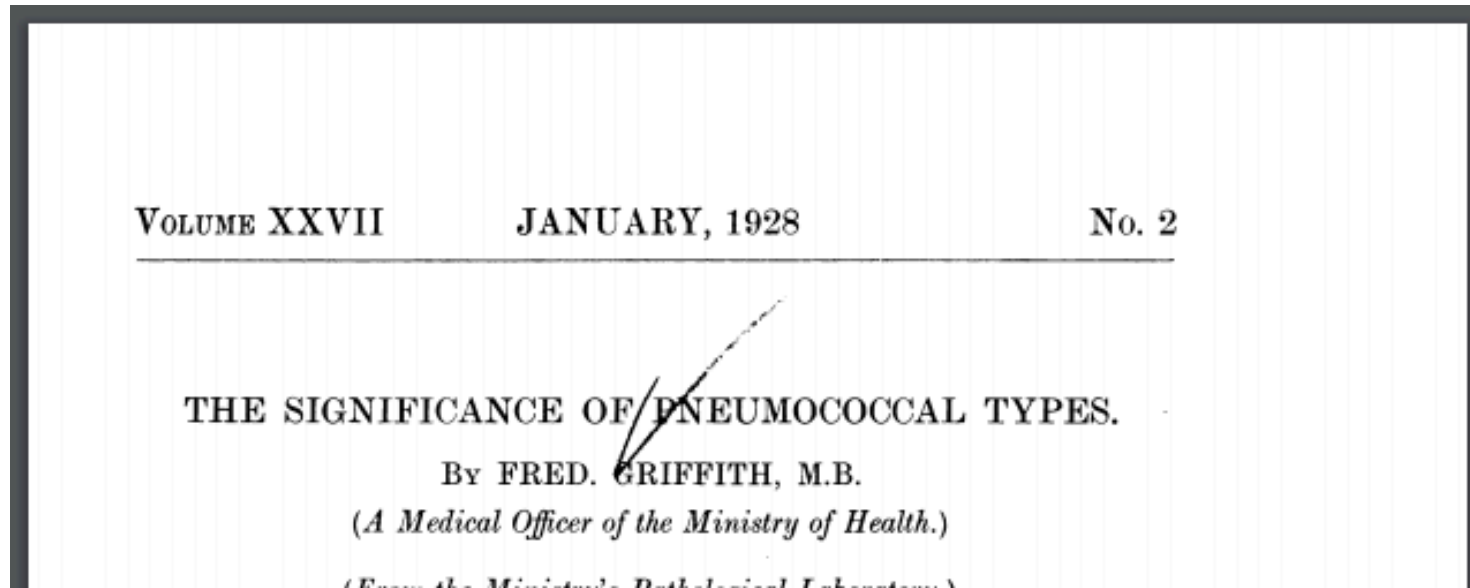
Dahm 2005

- Friedrich Miescher: isolated unknown substance from leukocyte nuclei
- Contained carbon, hydrogen, oxygen, nitrogen...
  - and high amounts of phosphorus, but no sulfur (so, it wasn't protein)
  - Called it "nuclein"

# DNA as the hereditary material

1928: Frederick Griffith shows that dead virulent bacteria can transform living non-virulent bacteria, making it virulent

What is the hereditary material that allows this?



# DNA as the hereditary material

1944: Oswald Avery isolates many different substances from virulent bacteria and applies them to nonvirulent bacteria

“Preparation 44” transforms the nonvirulent bacteria.

It is high in phosphorus, and more tests show it is DNA.

STUDIES ON THE CHEMICAL NATURE OF THE SUBSTANCE  
INDUCING TRANSFORMATION OF PNEUMOCOCCAL TYPES

INDUCTION OF TRANSFORMATION BY A DESOXYRIBONUCLEIC ACID FRACTION  
ISOLATED FROM PNEUMOCOCCUS TYPE III

By OSWALD T. AVERY, M.D., COLIN M. MacLEOD, M.D., AND  
MACLYN McCARTY,\* M.D.

*(From the Hospital of The Rockefeller Institute for Medical Research)*

# 1953: the Double Helix

737

## MOLECULAR STRUCTURE OF NUCLEIC ACIDS

### A Structure for Deoxyribose Nucleic Acid

WE wish to suggest a structure for the salt of deoxyribose nucleic acid (D.N.A.). This structure has novel features which are of considerable biological interest.

A structure for nucleic acid has already been proposed by Pauling and Corey<sup>1</sup>. They kindly made their manuscript available to us in advance of publication. Their model consists of three intertwined chains, with the phosphates near the fibre axis, and the bases on the outside. In our opinion, this structure is unsatisfactory for two reasons: (1) We believe that the material which gives the X-ray diagrams is the salt, not the free acid. Without the acidic hydrogen atoms it is not clear what forces would hold the structure together, especially as the negatively charged phosphates near the axis will repel each other. (2) Some of the van der Waals distances appear to be too small.

Another three-chain structure has also been suggested by Fraser (in the press). In his model the phosphates are on the outside and the bases on the inside, linked together by hydrogen bonds. This structure as described is rather ill-defined, and for this reason we shall not comment on it.

We wish to put forward a radically different structure for the salt of deoxyribose nucleic acid. This structure has two helical chains each coiled round the same axis (see diagram). We have made the usual chemical assumptions, namely, that each chain consists of phosphate diester groups joining  $\beta$ -D-deoxyribofuranose residues with 3',5' linkages. The two chains (but not their bases) are related by a dyad perpendicular to the fibre axis. Both chains follow right-handed helices, but owing to the dyad the sequences of the atoms in the two chains run in opposite directions. Each chain loosely resembles Furberg's<sup>2</sup> model No. 1; that is, the bases are on the inside of the helix and the phosphates on the outside. The configuration of the sugar and the atoms near it is close to Furberg's 'standard configuration', the sugar being roughly perpendicular to the attached base. There



This figure is purely diagrammatic. The two ribbons symbolize the two phosphate-sugar chains, and the horizontal rods the pairs of bases holding the chains together. The vertical line marks the fibre axis.

is a residue on each chain every 3.4 Å. in the z-direction. We have assumed an angle of 36° between adjacent residues in the same chain, so that the structure repeats after 10 residues on each chain, that is, after 34 Å. The distance of a phosphorus atom from the fibre axis is 10 Å. As the phosphates are on the outside, cations have easy access to them.

The structure is an open one, and its water content is rather high. At lower water contents we would expect the bases to tilt so that the structure could become more compact.

The novel feature of the structure is the manner in which the two chains are held together by the purine and pyrimidine bases. The planes of the bases are perpendicular to the fibre axis. They are joined together in pairs, a single base from one chain being hydrogen-bonded to a single base from the other chain, so that the two lie side by side with identical z-co-ordinates. One of the pair must be a purine and the other a pyrimidine for bonding to occur. The hydrogen bonds are made as follows: purine position 1 to pyrimidine position 1; purine position 6 to pyrimidine position 6.

If it is assumed that the bases only occur in the structure in the most plausible tautomeric forms (that is, with the keto rather than the enol configurations) it is found that only specific pairs of bases can bond together. These pairs are: adenine (purine) with thymine (pyrimidine), and guanine (purine) with cytosine (pyrimidine).

In other words, if an adenine forms one member of a pair, on either chain, then on these assumptions the other member must be thymine; similarly for guanine and cytosine. The sequence of bases on a single chain does not appear to be restricted in any way. However, if only specific pairs of bases can be formed, it follows that if the sequence of bases on one chain is given, then the sequence on the other chain is automatically determined.

It has been found experimentally<sup>3,4</sup> that the ratio of the amounts of adenine to thymine, and the ratio of guanine to cytosine, are always very close to unity for deoxyribose nucleic acid.

It is probably impossible to build this structure with a ribose sugar in place of the deoxyribose, as the extra oxygen atom would make too close a van der Waals contact.

The previously published X-ray data<sup>5,6</sup> on deoxyribose nucleic acid are insufficient for a rigorous test of our structure. So far as we can tell, it is roughly compatible with the experimental data, but it must be regarded as unproved until it has been checked against more exact results. Some of these are given in the following communications. We were not aware of the details of the results presented there when we devised our structure, which rests mainly though not entirely on published experimental data and stereochemical arguments.

It has not escaped our notice that the specific pairing we have postulated immediately suggests a possible copying mechanism for the genetic material. Full details of the structure, including the conditions assumed in building it, together with a set of co-ordinates for the atoms, will be published elsewhere.

We are much indebted to Dr. Jerry Donohue for constant advice and criticism, especially on interatomic distances. We have also been stimulated by a knowledge of the general nature of the unpublished experimental results and ideas of Dr. M. E. F. Wilkins, Dr. R. E. Franklin and their co-workers at

738

NATURE

April 25, 1953 VOL. 171

King's College, London. One of us (J. D. W.) has been aided by a fellowship from the National Foundation for Infantile Paralysis.

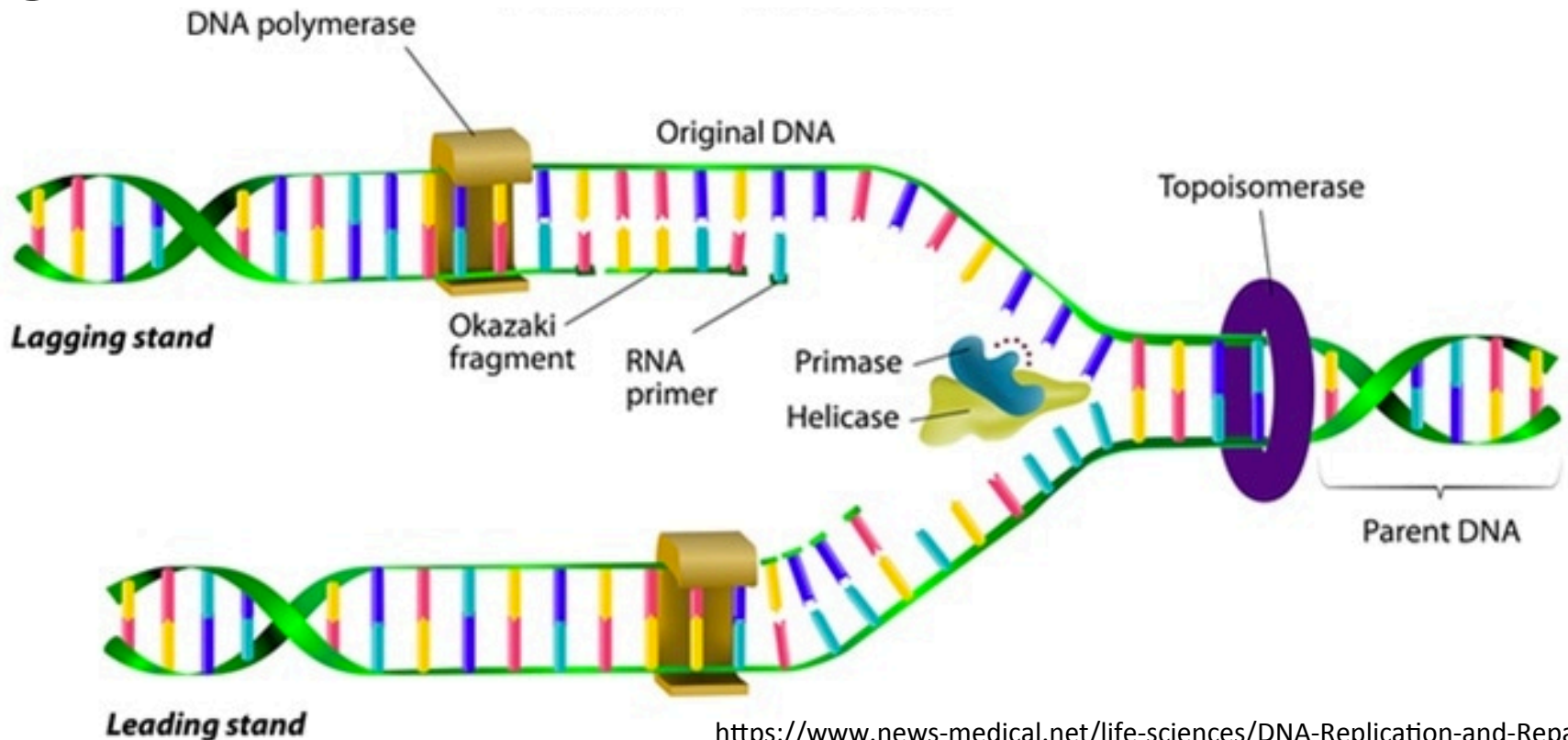
J. D. WATSON  
F. H. C. CRICK

Medical Research Council Unit for the Study of the Molecular Structure of Biological Systems, Cavendish Laboratory, Cambridge, April 2.

- <sup>1</sup> Pauling, L., and Corey, R. B. *Nature*, 171, 346 (1953); *Proc. U.S. Nat. Acad. Sci.*, 39, 54 (1953).  
<sup>2</sup> Furberg, S., *Acta Chem. Scand.*, 6, 634 (1952).  
<sup>3</sup> Chargaff, E., for references see Zamechok, S., Braverman, G., and Chargaff, E., *Biochim. et Biophys. Acta*, 9, 402 (1952).  
<sup>4</sup> Wyatt, G. R., *J. Gen. Physiol.*, 26, 201 (1952).  
<sup>5</sup> Astbury, W. T., *Symp. Soc. Exp. Biol.*, 1, *Nucleic Acid*, 66 (Camb. Univ. Press, 1947).  
<sup>6</sup> Wilkins, M. H. F., and Randall, J. T., *Biochim. et Biophys. Acta*, 10, 192 (1953).

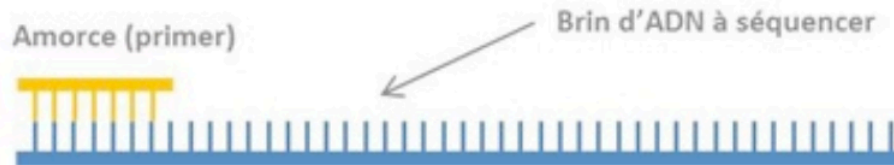
# DNA replication

“It has not escaped our notice that the specific pairing we have postulated immediately suggests a possible copying mechanism for the genetic material.”

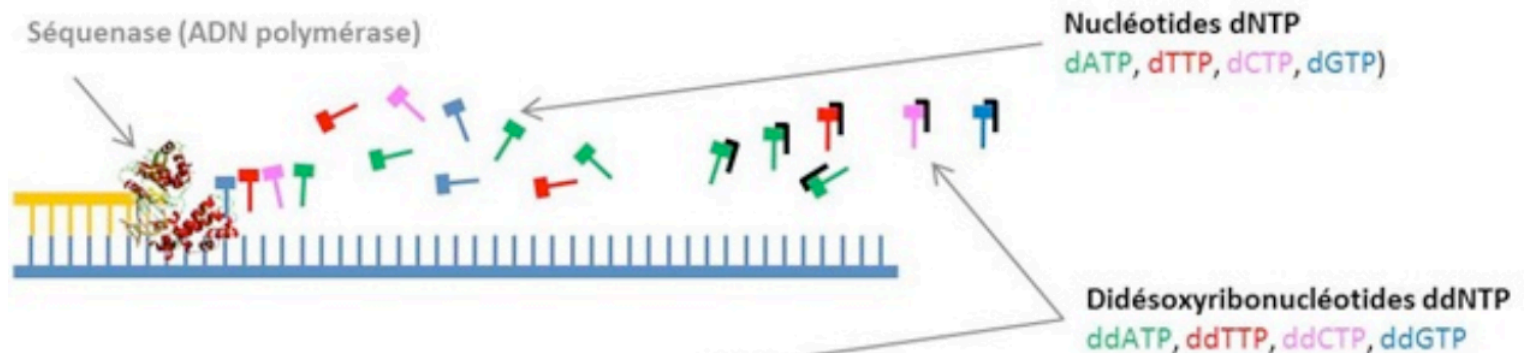


# 1977: Sanger's Chain-Termination method

Denaturation  
& priming



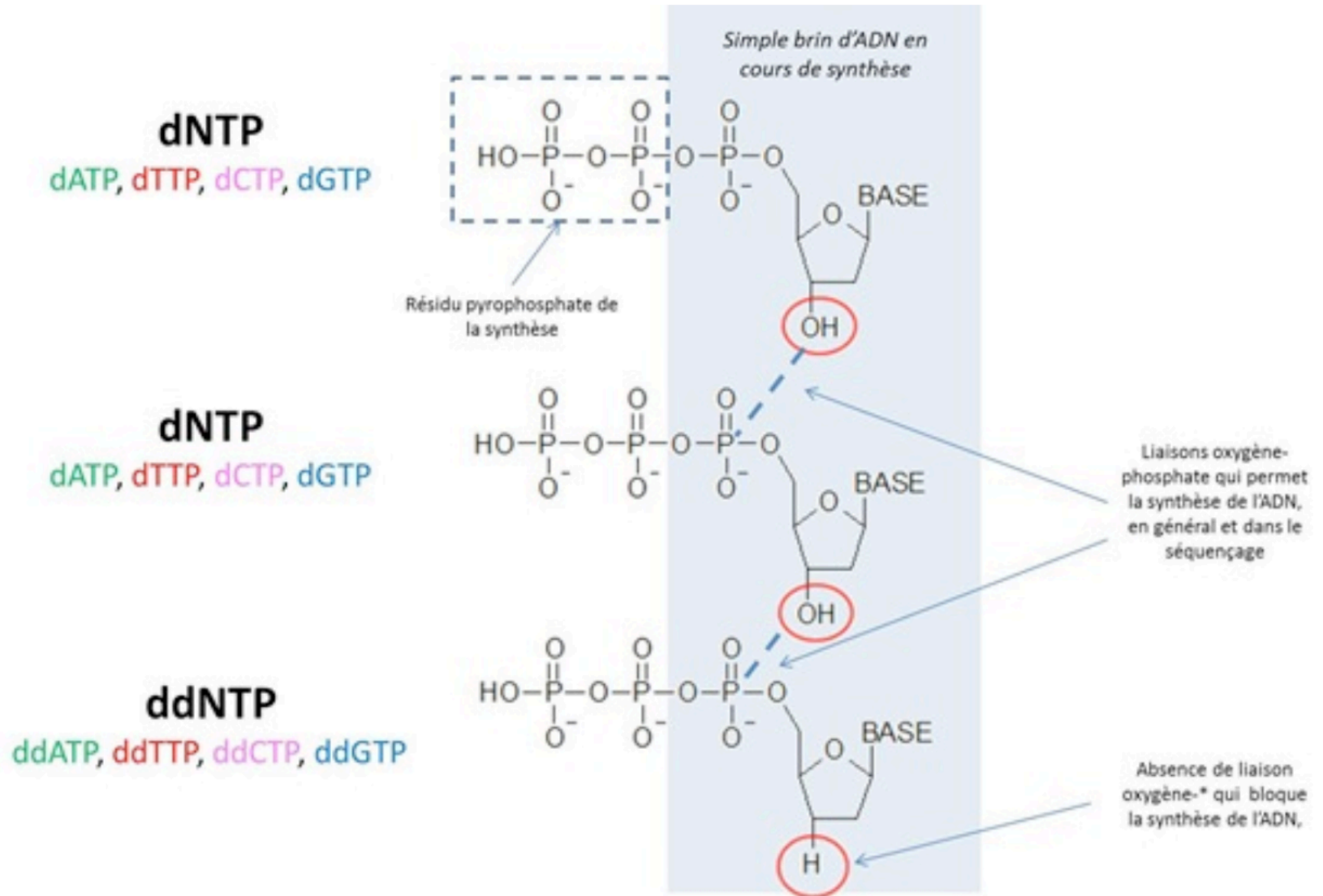
Polymerization



Arrest

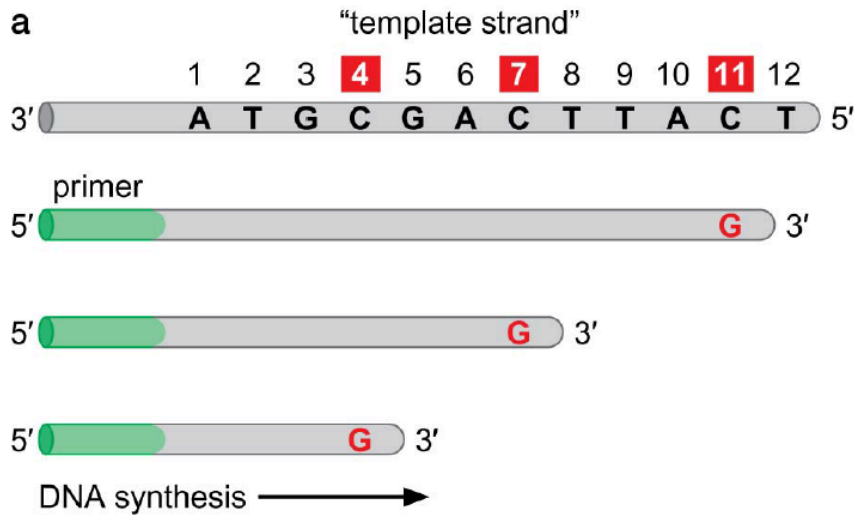


# 1977: Sanger's Chain-Termination method





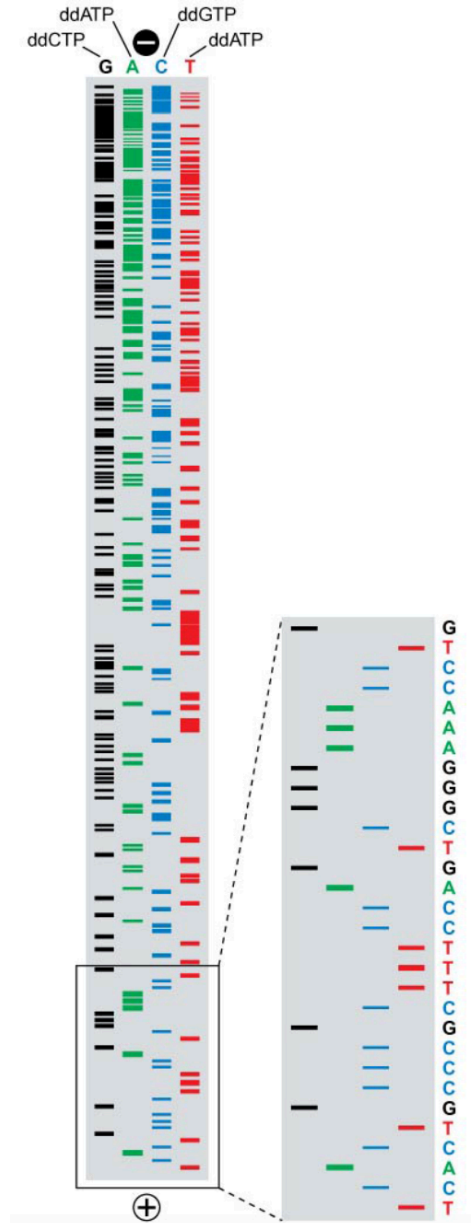
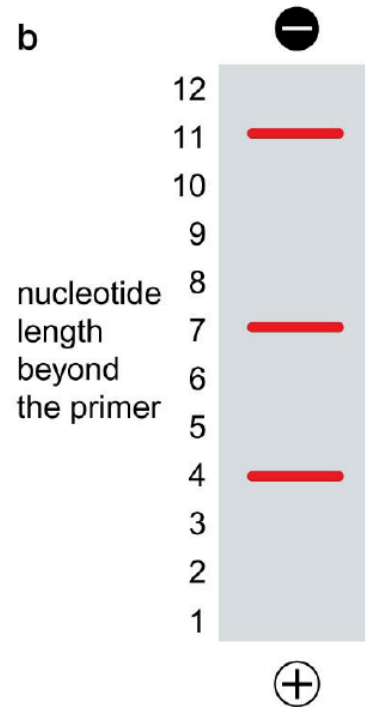
# 1977: Sanger's Chain-Termination method



substrates:

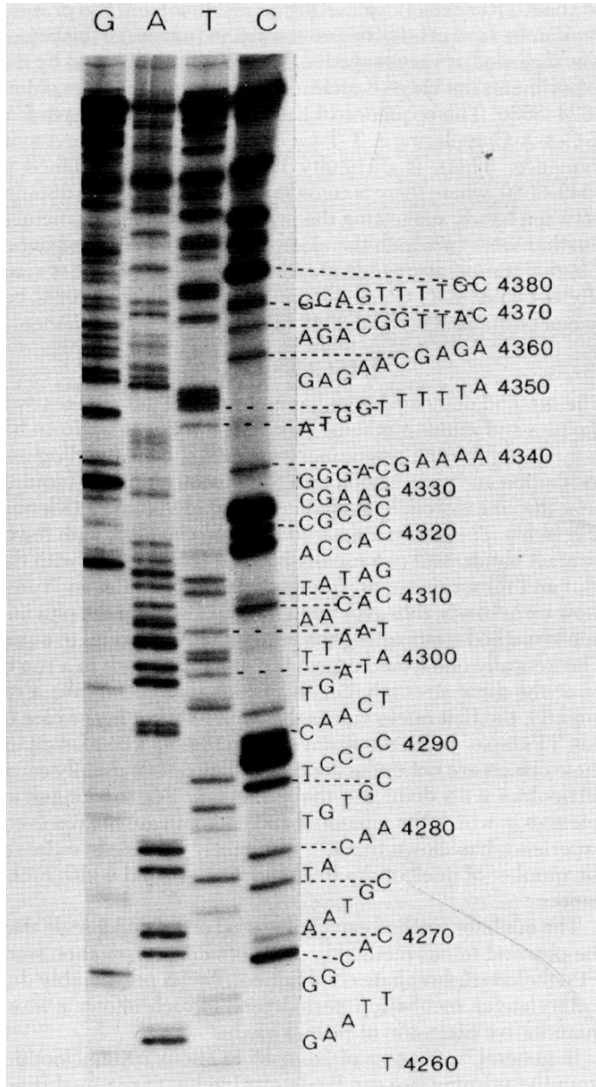
dATP dGTP ddGTP

dCTP dTTP



"sequencing by synthesis"

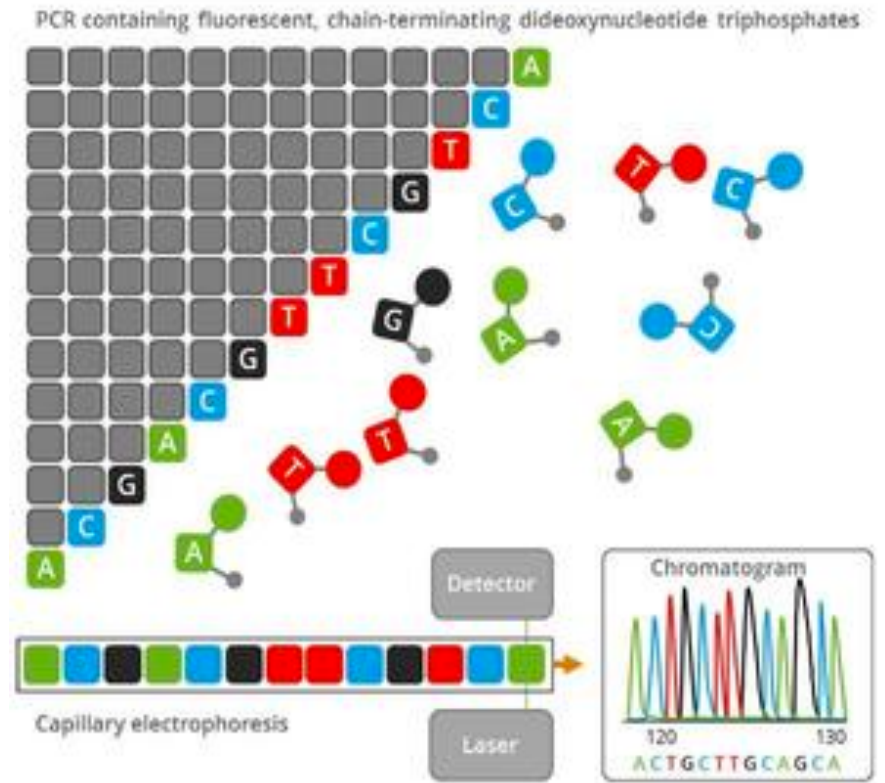
# 1977: Sanger's Chain-Termination method



“The electrophoresis was on a 12% acrylamide gel at 40 mA for 14 hr.”

# Improvements to the Sanger method

- Fluorophores instead of radio-labelling
- Capillary electrophoresis instead of acrylamide gels
- Discovery of more useful enzymes
- Automation



Format: Abstract ▾

Send to ▾

[Cell](#). 1993 Mar 26;72(6):971-83.**A novel gene containing a trinucleotide repeat that is expanded and unstable on Huntington's disease chromosomes. The Huntington's Disease Collaborative Research Group.**

[No authors listed]

**Abstract**

The Huntington's disease (HD) haplotype analysis of linkage disequilibrium revealed a new gene, IT15, isolated using a trinucleotide repeat that is expanded and unstable on HD chromosomes and 4p16.3 haplotypes. The (CAG)<sub>n</sub> repeat encodes an approximately 348 kd protein that is involved in myotonic dystrophy, acting

**Comment in**

Planting alfalfa and cloning the

PMID: 8458085

[Indexed for MEDLINE]



## Identification of the Cystic Fibrosis Gene: Cloning and Characterization of the Complete Coding Sequence

JOHN R. RIORDAN, JOHANNA M. Riordan, RICHARD ROZMAHEL, ZBYSZEK R. ZBYSZEK, NATASA PLAVSIC, JIA-LING CHOU, FRANCIS S. COLLINS, LAP-CHEE TSUI

Overlapping complementary DNA clones were isolated from epithelial cell libraries with a genomic DNA segment containing a portion of the putative cystic fibrosis (CF) locus, which is on chromosome 7. Transcripts, approximately 6500 nucleotides in size, were detectable in the tissues affected in patients with CF. The predicted protein consists of two similar motifs, each with (i) a domain having properties consistent with membrane association and (ii) a domain believed to be involved in ATP (adenosine triphosphate) binding. A deletion of three

isolation of polypeptide complement that mediates conductance activated pathway and CF biochemical defect in CF remains to be determined.

Molecular cloning experiments revealed a large, contiguous segment of DNA containing the sequences from a region of approximately 7 kb (7). These sequences were in a region of DNA that was able to detect conserved sequences by Southern DNA hybridization and were confirmed by Northern hybridization experiments, cl

MENU ▾

**nature**  
International journal of science

Article | Published: 24 January 1985

## Complete nucleotide sequence of the AIDS virus, HTLV-III

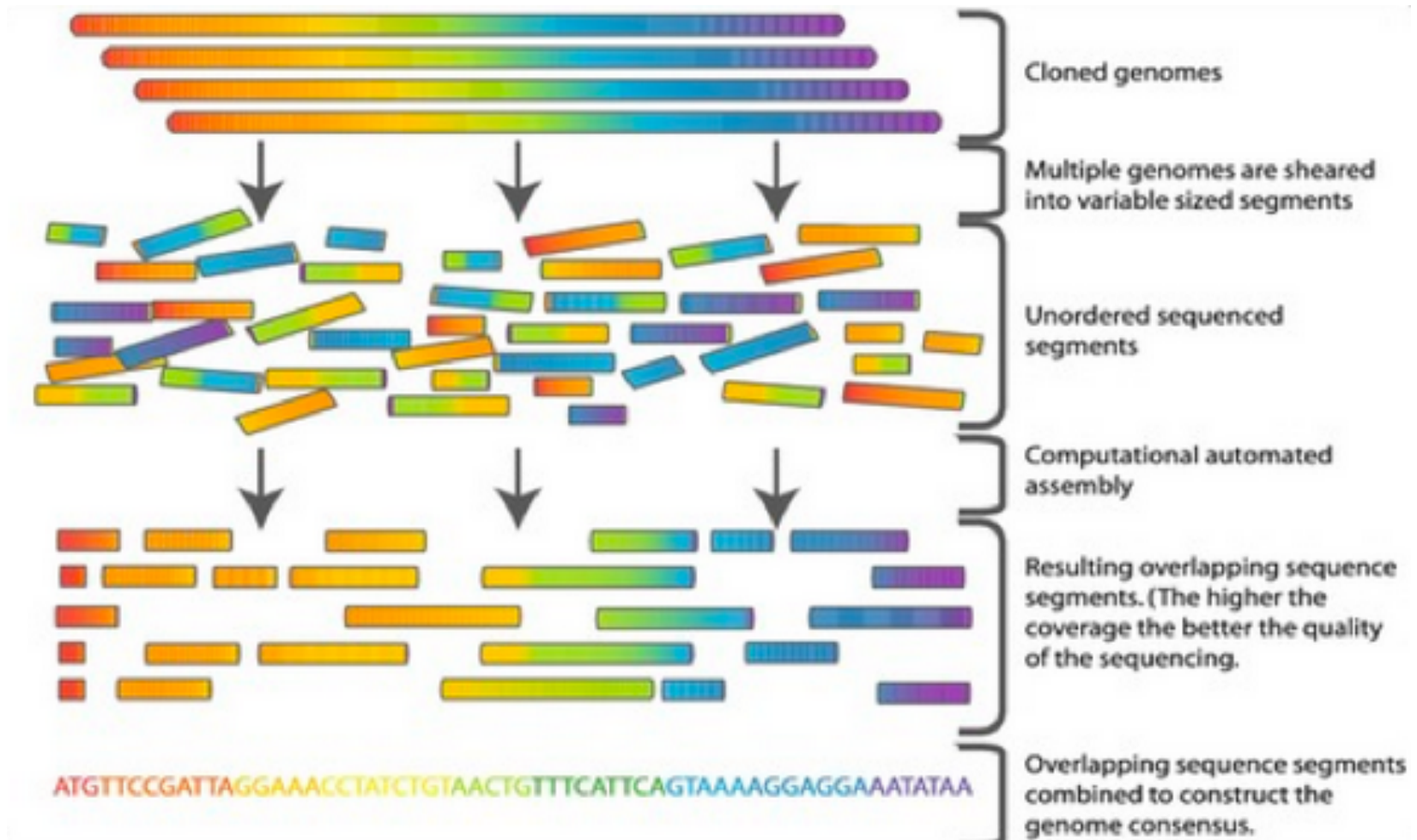
Lee Ratner, William Haseltine, Roberto Patarca, Kenneth J. Livak, Bruno Starcich, Steven F. Joseph, Ellen R. Doran, J. Antoni Rafalski, Erik A. Whitehorn, Kirk Baumeister, Lucinda Ivanoff, Stephen R. Petteway Jr, Mark L. Pearson, James A. Lautenberger, Takis S. Papas, John Ghrayeb, Nancy T. Chao, Robert C. Gallo & Flossie Wong-Staal

*Nature* **313**, 277–284 (24 January 1985) | [Download Citation](#) ↓

### Abstract

The complete nucleotide sequence of two human T-cell leukaemia type III (HTLV-III) proviral DNAs each have four long open reading frames, the first two corresponding to the gag and pol genes. The fourth open reading frame encodes two functional polypeptides, a large precursor of the major envelope glycoprotein and a smaller protein derived from the 3'-terminus long open reading frame analogous to the long open reading frame (*lor*) product of HTLV-I and -II.

# Shotgun Sequencing



# Human Genome Project

- Initiated in 1990
  - Goal: finished genome by 2005
- Hundreds of labs in 18 countries
  - Sequences made available immediately after assembly
- Competitors: J. Craig Venter and Celera

## articles

### Initial sequencing and analysis of the human genome

International Human Genome Sequencing Consortium\*

\* A partial list of authors appears on the opposite page. Affiliations are listed on the next page.

The human genome holds an extraordinary trove of information. Here we report the results of an international collaboration to sequence and assemble a high-quality draft of the human genome. We also present an initial analysis of the data, describing the structure and organization of the genome.

The rediscovery of Mendel's laws of heredity in the opening weeks of the 20th century<sup>1-3</sup> sparked a scientific quest to understand the nature and content of genetic information that has propelled biology for the last hundred years. The scientific progress made falls naturally into four main phases, corresponding roughly to four quarters of the century. The first established the cellular basis of heredity: the chromosomes. The second defined the molecular basis of heredity: the DNA double helix. The third unlocked the information

## THE HUMAN GENOME

### The Sequence of the Human Genome

## articles

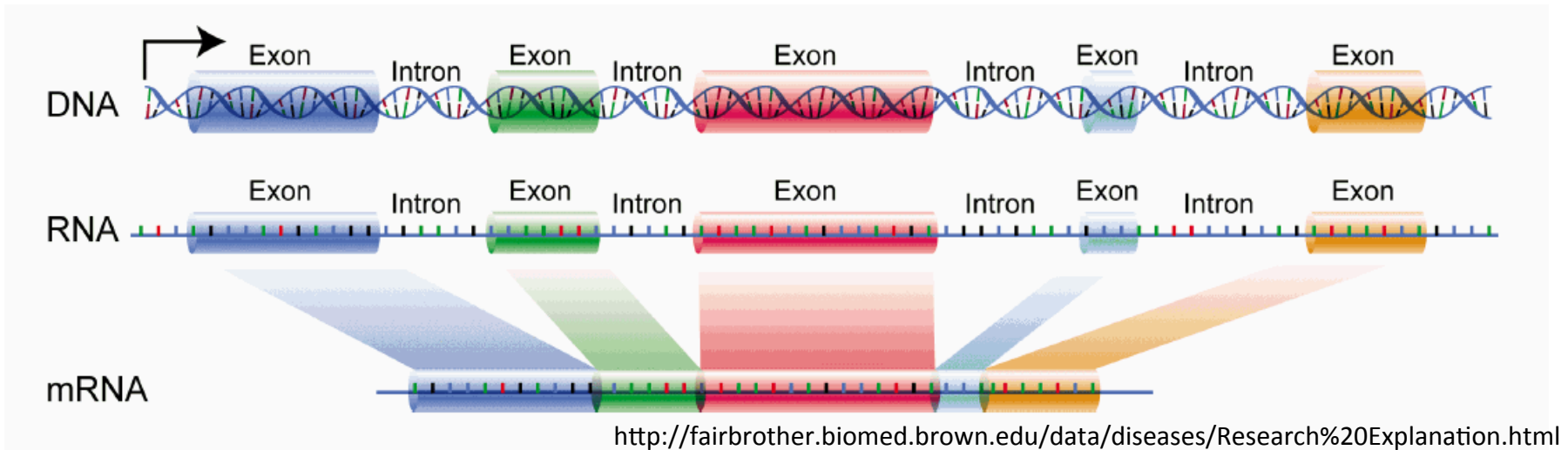
### Finishing the euchromatic sequence of the human genome

International Human Genome Sequencing Consortium\*

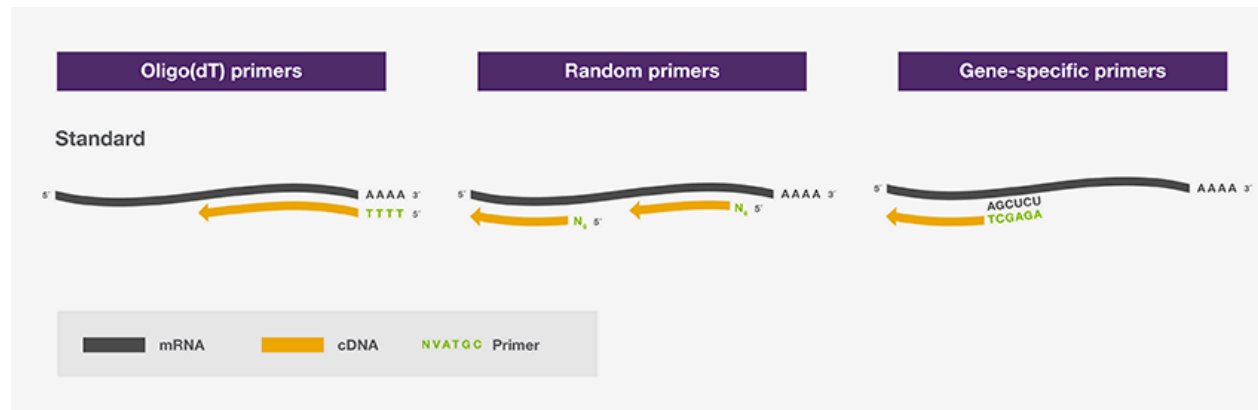
\* A list of authors and their affiliations appears in the Supplementary Information

Gene W. Myers,<sup>1</sup> Peter W. Li,<sup>1</sup> Richard J. Mural,<sup>1</sup> Mark Yandell,<sup>1</sup> Cheryl A. Evans,<sup>1</sup> Robert A. Holt,<sup>1</sup> Michael S. Waterman,<sup>1</sup> Richard M. Ballew,<sup>1</sup> Daniel H. Huson,<sup>1</sup> John L. M. Dunham,<sup>1</sup> Shinnappa D. Kodira,<sup>1</sup> Xiangqun H. Zheng,<sup>1</sup> Lin Chen,<sup>1</sup> Pradyumn K. Jain,<sup>1</sup> Paul D. Thomas,<sup>1</sup> Jinghui Zhang,<sup>1</sup> Jianzhong Zhang,<sup>1</sup> Steven A. Wheeler,<sup>1</sup> Samuel Broder,<sup>1</sup> Andrew G. Clark,<sup>4</sup> Joe Nadeau,<sup>5</sup> Arnold J. Levine,<sup>7</sup> Richard J. Roberts,<sup>8</sup> Mel Simon,<sup>9</sup> David G. Klapper,<sup>1</sup> Richard K. Wilson,<sup>1</sup> Randall Bolanos,<sup>1</sup> Arthur Delcher,<sup>1</sup> Ian Dew,<sup>1</sup> Daniel Fasulo,<sup>1</sup> Glenn S. Gager,<sup>1</sup> Sridhar Hannenhalli,<sup>1</sup> Saul Kravitz,<sup>1</sup> Samuel Levy,<sup>1</sup> Jerome A. McClellan,<sup>1</sup> Jane Abu-Threideh,<sup>1</sup> Ellen Beasley,<sup>1</sup> Kendra Biddick,<sup>1</sup> James C. Binkley,<sup>1</sup> Richard W. Burch,<sup>1</sup> David W. Cargill,<sup>1</sup> Ishwar Chandramouliswaran,<sup>1</sup> Rosane Charlab,<sup>1</sup> Stephen Chen,<sup>1</sup> Mark A. Chertney,<sup>1</sup> Antina Di Francesco,<sup>1</sup> Patrick Dunn,<sup>1</sup> Karen Eilbeck,<sup>1</sup> David R. Fulton,<sup>1</sup> Jianjun Gan,<sup>1</sup> Wangmao Ge,<sup>1</sup> Fangcheng Gong,<sup>1</sup> Zhiping Gu,<sup>1</sup> Robert Guiguen, <sup>1</sup> James E. Higgins,<sup>1</sup> Rui-Ru Ji,<sup>1</sup> Zhaoxi Ke,<sup>1</sup> Karen A. Ketchum,<sup>1</sup> Richard A. Kittling,<sup>1</sup> John D. Kim,<sup>1</sup> Jianjun Li,<sup>1</sup> Yong Liang,<sup>1</sup> Xiaoying Lin,<sup>1</sup> Fu Lu,<sup>1</sup>

# Sequencing RNA

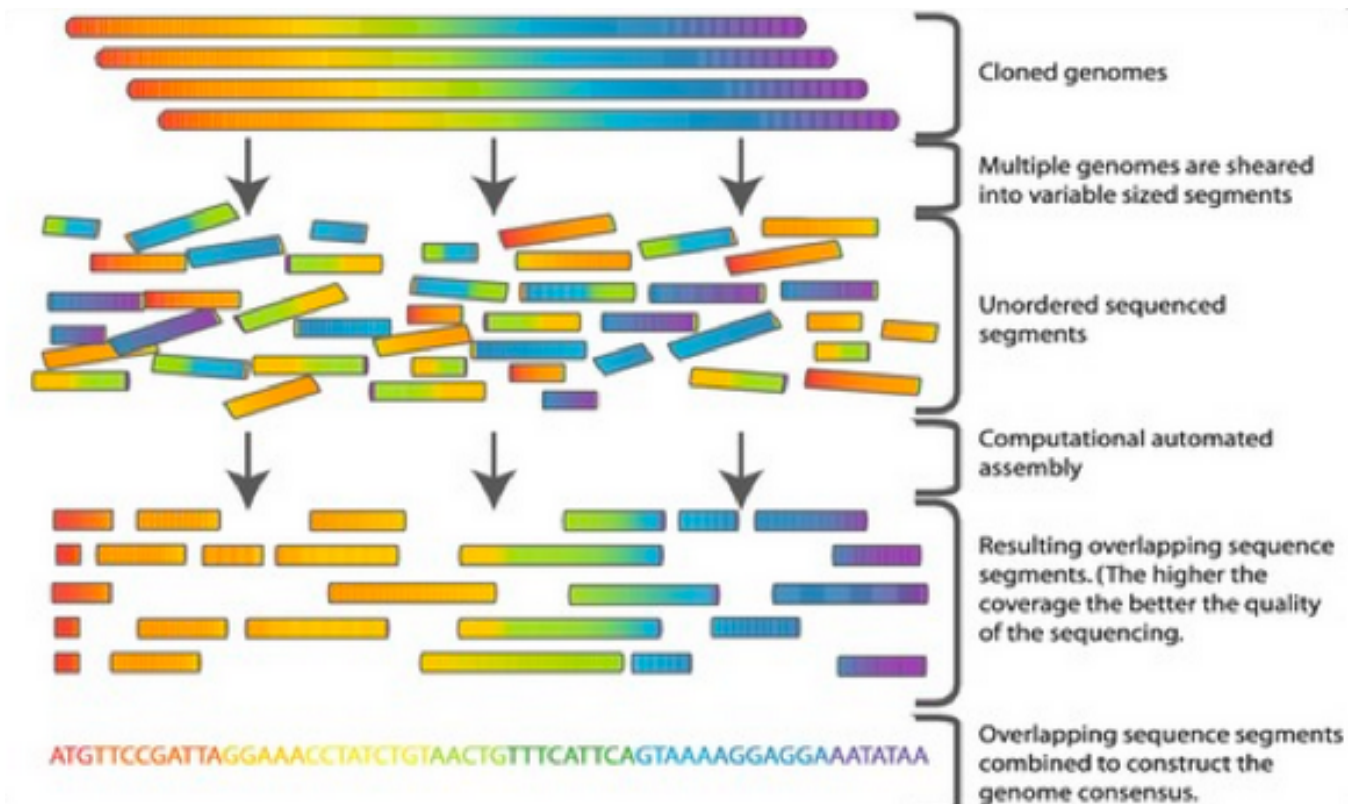


Reverse transcription to cDNA is (almost) always necessary.



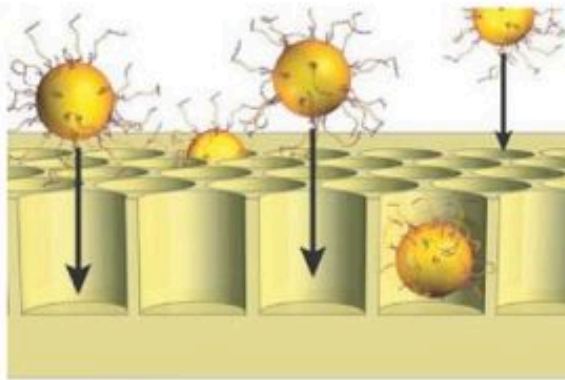
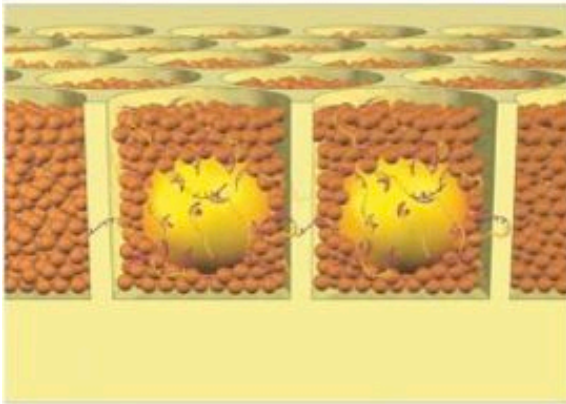
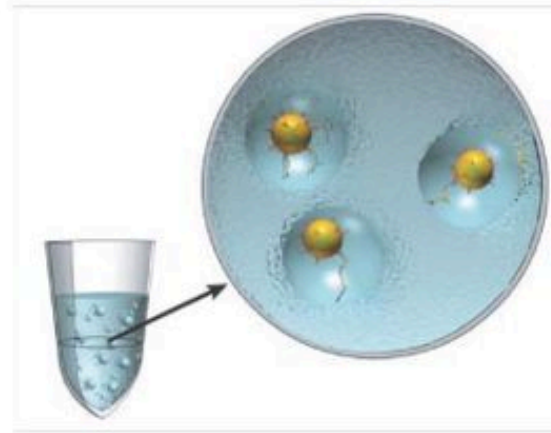
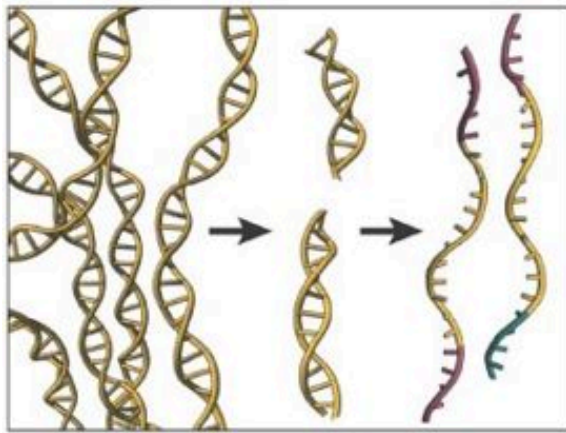
# Second Generation Sequencing

- AKA “next generation sequencing” or NGS
- Massively parallel sequencing
- Short read length

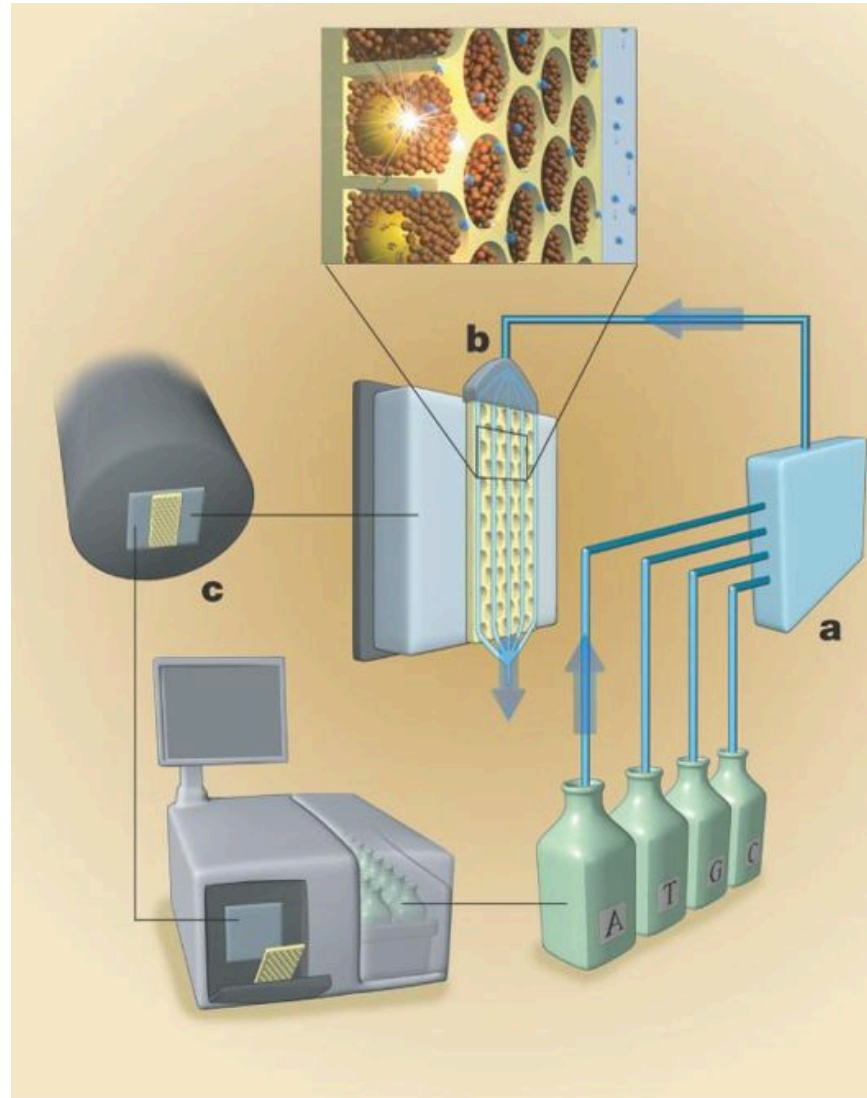




# ~2005: 454 Pyrosequencing



# 454 Pyrosequencing



# Genome Sequencing with 454

“Here we report the DNA sequence of a diploid genome of a single individual, James D. Watson, sequenced to 7.4-fold redundancy in **two months** using massively parallel sequencing in picolitre-size reaction vessels.”

nature

Vol 452 | 17 April 2008 | doi:10.1038/nature06884

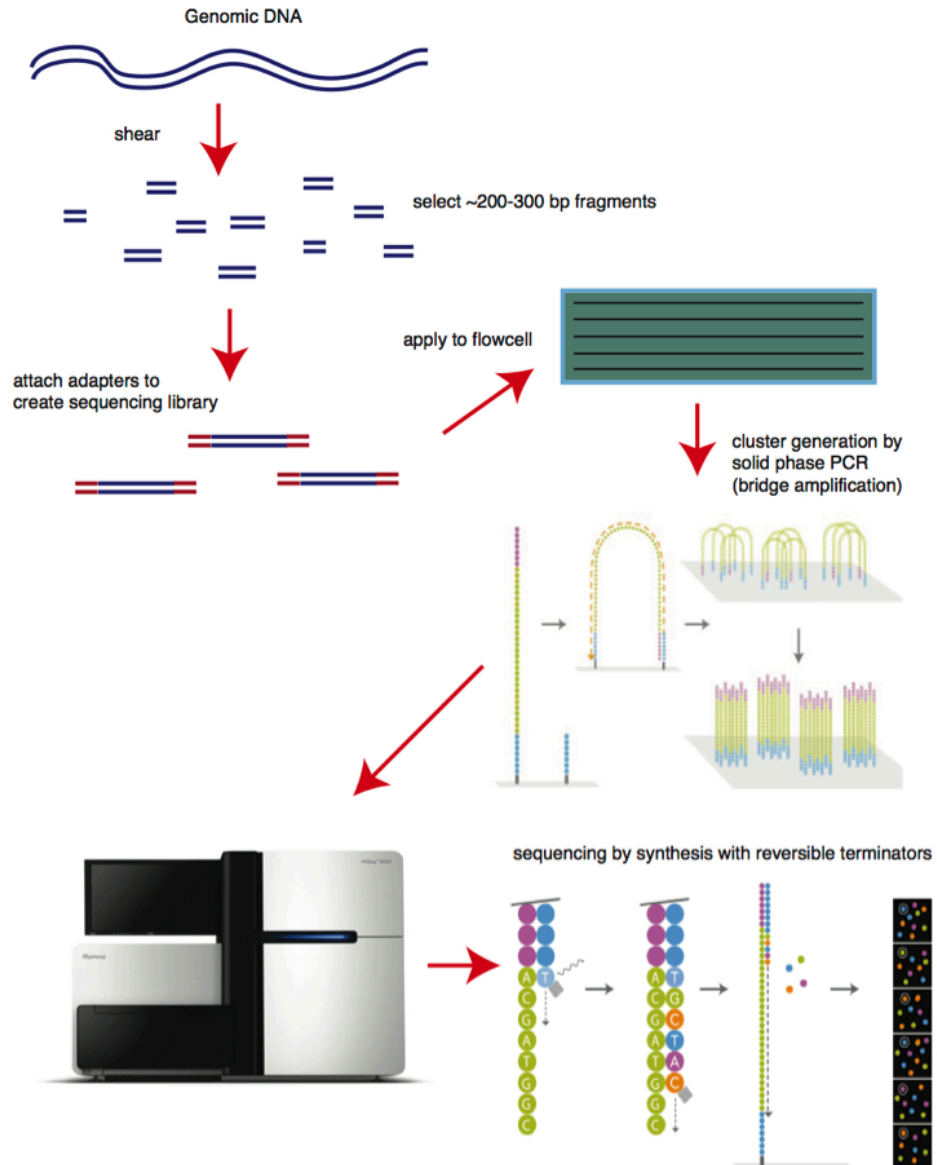
## LETTERS

---

### **The complete genome of an individual by massively parallel DNA sequencing**

David A. Wheeler<sup>1\*</sup>, Maithreyan Srinivasan<sup>2\*</sup>, Michael Egholm<sup>2\*</sup>, Yufeng Shen<sup>1\*</sup>, Lei Chen<sup>1</sup>, Amy McGuire<sup>3</sup>, Wen He<sup>2</sup>, Yi-Ju Chen<sup>2</sup>, Vinod Makhiani<sup>2</sup>, G. Thomas Roth<sup>2</sup>, Xavier Gomes<sup>2</sup>, Karrie Tartaro<sup>2†</sup>, Faheem Niazi<sup>2</sup>.

# ~2006: Solexa/Illumina



# An integrated map of genetic variation from 1,092 human genomes

The 1000 Genomes Project Consortium\*

By characterizing the geographic and functional spectrum of human genetic variation, the 1000 Genomes Project aims to build a resource to help to understand the genetic contribution to disease. Here we describe the genomes of 1,092 individuals from 14 populations, constructed using a combination of low-coverage whole-genome and exome sequencing. By developing methods to integrate information across several algorithms and diverse data sources, we provide a validated haplotype map of 38 million single nucleotide polymorphisms, 1.4 million short insertions and deletions, and more than 14,000 larger deletions. We show that individuals from different populations carry different profiles of rare and common variants, and that low-frequency variants show substantial geographic differentiation, which is further increased by the action of purifying selection. We show that evolutionary conservation and coding consequence are key determinants of the strength of purifying selection, that rare-variant load varies substantially across biological pathways, and that each individual contains hundreds of rare variants such as motif-disrupting changes in transcription-factor-binding sites. This accessible single nucleotide polymorphisms at a frequency of 1% in related populations and low-frequency variants in individuals from diverse, including admixed, populations

Recent efforts to map human genetic variation by sequencing exomes<sup>1</sup> and whole genomes<sup>2-4</sup> have characterized the vast majority of common single nucleotide polymorphisms (SNPs) and many structural variants across the genome. However, although more than 95% of common (>5% frequency) variants were discovered in the pilot phase of the 1000 Genomes Project, lower-frequency variants, particularly those outside the coding exome, remain poorly characterized. Low-frequency variants are enriched for potentially functional mutations, for example, protein-changing variants, under weak purifying selection<sup>5,6</sup>. Furthermore, because low-frequency variants tend to be recent in origin, they exhibit increased levels of population differentiation<sup>6-8</sup>. Characterizing such variants, for both point mutations and structural changes, across a range of populations is thus likely to identify many variants of functional importance and is crucial for interpreting

individual genome those private to fact

We now report on the genomes of 1,092 individuals sampled from 14 populations drawn from Europe, East Asia, sub-Saharan Africa and the Americas (Supplementary Figs 1 and 2), analysed through a combination of low-coverage (2-6×) whole-genome sequence data, targeted deep (50-100×) exome sequence data and dense SNP genotype data (Table 1 and Supplementary Tables 1-3). This design was shown by the pilot phase<sup>2</sup> to be powerful and cost-effective in discovering and genotyping all but the rarest SNP and short insertion and deletion (indel) variants. Here, the approach was augmented with statistical methods for selecting higher quality variant calls from candidates obtained using multiple algorithms, and to integrate SNP, indel and larger structural variants within a single framework (see

**Table 1 | Summary of 1000 Genomes Project phase I data**

	Autosomes	Chromosome X	GENCODE regions*
Samples	1,092	1,092	1,092
Total raw bases (Gb)	19,049	804	327
Mean mapped depth (×)	5.1	3.9	80.3
SNPs			
No. sites overall	36.7 M	1.3 M	498 K
Novelty rate†	58%	77%	50%
No. synonymous/non-synonymous/nonsense	NA	4.7/6.5/0.097 K	199/293/6.3 K
Average no. SNPs per sample	3.60 M	105 K	24.0 K
Indels			
No. sites overall	1.38 M	59 K	1,867
Novelty rate†	62%	73%	54%
No. in-frame/frameshift	NA	19/14	719/1,066
Average no. indels per sample	344 K	13 K	440
Genotyped large deletions			
No. sites overall	13.8 K	432	847
Novelty rate†	54%	54%	50%
Average no. variants per sample	717	26	39

NA, not applicable.

\*Autosomal genes only.

†Compared with dbSNP release 135 (Oct 2011), excluding contribution from phase I 1000 Genomes Project (or equivalent data for large deletions).

\*Lists of participants and their affiliations appear at the end of the paper.



## GENETIC RESEARCH

### The 100 000 Genomes Project

Part research project, part commercial stimulus, this enormous sequencing programme could usher genomic medicine into mainstream use, **Mark Peplow** reports

Mark Peplow *freelance journalist, Cambridge*

## FEATURE

# Variant analysis

is now ramping up into high gear. Overseen by Genomics England, it is one of the biggest whole genome sequencing projects in the world. And it is working to a breathtaking timetable: most of these genomes will be sequenced by the end of next year.

The genetic material will come from patients with rare diseases or common cancers and their families (box 1). By identifying any genetic anomalies, and linking them to participants' medical histories for the rest of their lives, the project aims to build up a unique database for treatment and research. "It will allow us to find things in the data that we might not notice in ordinary clinical care," says Caulfield. That should offer better diagnoses and more targeted therapies. It also gives scientists a treasure trove of information that could help to develop more effective drugs.

That remit is impressive enough. But the project's broader goals are to kickstart a national genomics industry and make the UK the first country to routinely use DNA sequencing in mainstream healthcare. "If we get this right, our ambition is to see new treatments, new diagnostics, coming to patients in the UK first," says Caulfield.

### Clinical potential

The project is already having clinical impact among people with rare diseases, with the first child participants receiving a genetic diagnosis in January. There are about 7000 known rare diseases, and roughly 1 in 17 people (about three million in the UK) are affected at some point in their lives.<sup>1</sup> "Collectively, the burden is high," says Caulfield. "They are a huge cause of disability, and the toll on individuals is huge."

More than 80% of rare diseases are suspected to have a genetic component. But their rarity makes diagnosis a huge challenge,

Beverly Searle, chief executive of Unique, the rare chromosome disorder support group ([www.rarechromo.co.uk](http://www.rarechromo.co.uk)). "That's why it's so important to have these projects where you gather large datasets." With tens of thousands of genomes from patients with rare diseases, it becomes much more likely to find a statistically robust association between genetic variants and a particular disease.

Once that link is established, the experiences of those who share the same genetic anomalies can be compared to predict how a particular patient's condition might develop in the future and which treatments are likely to be more effective. Not only could this improve clinical outcomes, it could also save time and money. "It's important that expectations are managed—not every family will get a diagnosis," cautions Searle. "But we've got the potential for one test to give you an answer."

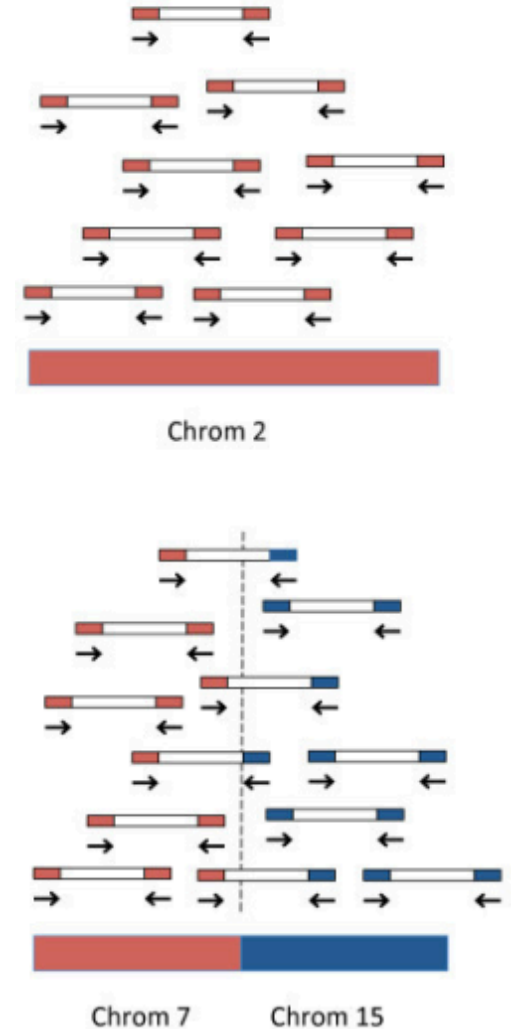
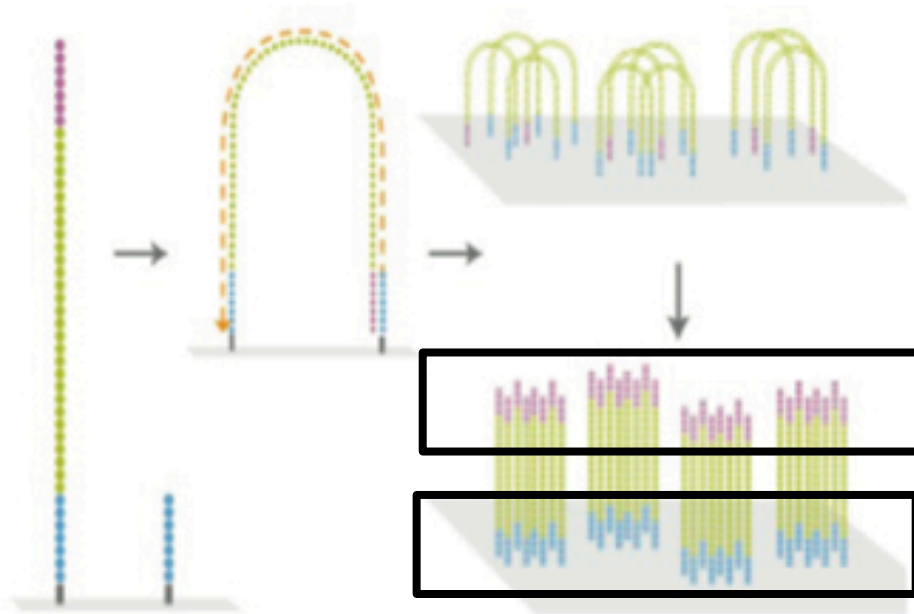
The other arm of the project focuses on common cancers, including those in the lung, breast, colon, prostate, and ovary, where a genetic diagnosis could affect treatment options. About half of melanomas are caused by a mutation in the *BRAF* gene, for example, and these can be treated with a drug that specifically targets the BRAF protein. "A mutation can help to predict a medicine's effectiveness," says Caulfield.

### Collecting and using the data

Most of the project's participants arrive via one of the 13 NHS Genomic Medicine Centres that were established last year around England. People give a small blood sample, and (if they have cancer) a small piece of their tumour, which can have a substantially different genome.

The project has already sequenced more than 7000 genomes and is recruiting more than 200 patients with rare diseases per week. But the pace of sequencing will quicken in the coming months, when a dedicated facility at the Wellcome Genome Campus in Hinxton, Cambridgeshire, opens. The American company Illumina is setting up a world class sequencing facility there, stuffed with machines that can read an entire genome in

# Paired-end sequencing



<https://bitesizebio.com/13546/sequencing-by-synthesis-explaining-the-illumina-sequencing-technology/>

## ARTICLES

# A human gut microbial gene catalogue established by metagenomic sequencing

Junjie Qin<sup>1\*</sup>, Ruiqiang Li<sup>1\*</sup>, Jeroen Raes<sup>2,3</sup>, Manimozhayan Arumugam<sup>2</sup>, Kristoffer Solvsten Burgdorf<sup>4</sup>, Chaysavanh Manichanh<sup>5</sup>, Trine Nielsen<sup>4</sup>, Nicolas Pons<sup>6</sup>, Florence Levenez<sup>6</sup>, Takuji Yamada<sup>2</sup>, Daniel R. Mende<sup>7</sup>, Junhua Li<sup>1,7</sup>, Junming Xu<sup>1</sup>, Shaochuan Li<sup>1</sup>, Dongfang Li<sup>1,8</sup>, Jianjun Cao<sup>1</sup>, Bo Wang<sup>1</sup>, Huiqing Liang<sup>1</sup>, Huisong Zher<sup>1</sup>, Yinlong Xie<sup>1,7</sup>, Julien Tap<sup>6</sup>, Patricia Lepage<sup>6</sup>, Marcelo Bertalan<sup>9</sup>, Jean-Michel Batto<sup>6</sup>, Torben Hansen<sup>4</sup>, Denis L. Paslier<sup>10</sup>, Allan Linneberg<sup>11</sup>, H. Bjørn Nielsen<sup>9</sup>, Eric Pelletier<sup>10</sup>, Pierre Renault<sup>6</sup>, Thomas Sicheritz-Ponten<sup>9</sup>, Keith Turner<sup>12</sup>, Hongmei Zhu<sup>1</sup>, Chang Yu<sup>1</sup>, Shengting Li<sup>1</sup>, Min Jian<sup>1</sup>, Yan Zhou<sup>1</sup>, Yingrui Li<sup>1</sup>, Xiuqing Zhang<sup>1</sup>, Songgang Li<sup>1</sup>, Nan Qin<sup>1</sup>, Huanming Yang<sup>1</sup>, Jian Wang<sup>1</sup>, Søren Brunak<sup>9</sup>, Joel Doré<sup>6</sup>, Francisco Guarner<sup>5</sup>, Karsten Kristiansen<sup>13</sup>, Oluf Pedersen<sup>4,14</sup>, Julian Parkhill<sup>12</sup>, Jean Weissenbach<sup>10</sup>, MetaHIT Consortium†, Peer B. S. Dusko Ehrlich<sup>6</sup> & Jun Wang<sup>1,13</sup>

To understand the impact of gut microbes on human health we describe the Illumina-based metagenomic sequencing, microbial genes, derived from 576.7 gigabases of sequence ~150 times larger than the human gene complement, containing microbial genes of the cohort and probably includes a large number of genes that are largely shared among individuals of the cohort. The cohort harbours between 1,000 and 1,150 prevalent bacterial species and each individual at least 160 such species, which are also largely shared. We define and describe the minimal gut metagenome and the minimal gut bacterial genome in terms of functions present in all individuals and most bacteria, respectively.

It has been estimated that the microbes in our bodies collectively make up to 100 trillion cells, tenfold the number of human cells, and suggested that they encode 100-fold more unique genes than our own genome<sup>1</sup>. The majority of microbes reside in the gut, have a profound influence on human physiology and nutrition, and are crucial for human life<sup>2,3</sup>. Furthermore, the gut microbes contribute to energy harvest from food, and changes of gut microbiome may be associated with bowel diseases or obesity<sup>4-6</sup>.

To understand and exploit the impact of the gut microbes on human health and well-being it is necessary to decipher the content, diversity and functioning of the microbial gut community. 16S ribosomal RNA gene (rRNA) sequence-based methods<sup>7</sup> revealed that two bacterial divisions, the Bacteroidetes and the Firmicutes, constitute over 90% of the known phylogenetic categories and dominate the distal gut microbiota<sup>10</sup>. Studies also showed substantial diversity of the gut microbiome between healthy individuals<sup>4,8,10,11</sup>. Although this difference is especially marked among infants<sup>12</sup>, later in life the gut microbiome converges to more similar phyla.

Metagenomic sequencing represents a powerful alternative to rRNA sequencing for analysing complex microbial communities<sup>13-15</sup>. Applied to the human gut, such studies have already generated some 3 gigabases (Gb) of microbial sequence from faecal samples of 33

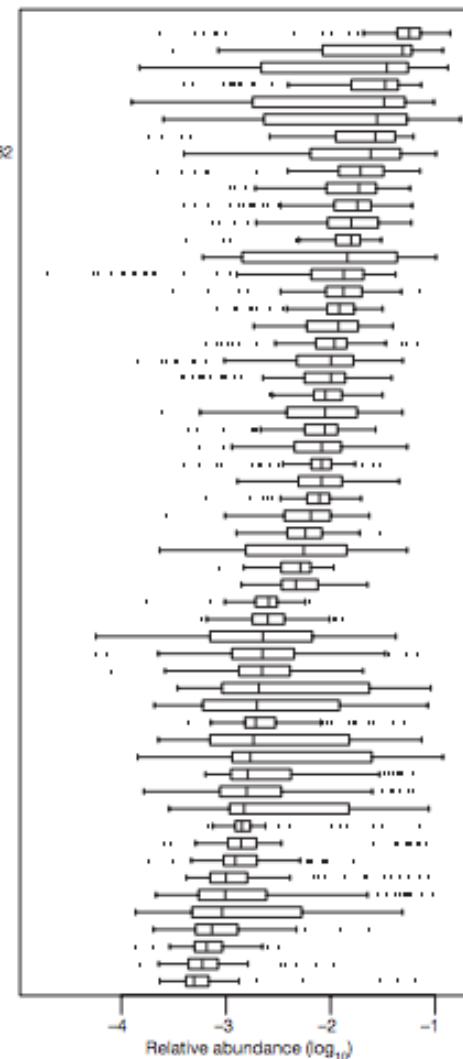
individuals from the United States or Japan<sup>8,16,17</sup>. To get a broad overview of the human gut microbial genes we used the Illumina Genome Analyser (GA) technology to carry out deep sequencing of total DNA from faecal samples of 124 European adults. We generated 576.7 Gb of sequence, almost 200 times more than in all previous studies, assembled it into contigs and predicted 3.3 million unique open reading frames (ORFs). This gene catalogue contains virtually all of the prevalent gut microbial genes in our cohort, providing a broad view of the functions important for bacterial life in the gut and indicates that many bacterial species are shared by different individuals. Our results also show that short-read metagenomic sequencing can be used for global characterization of the genetic potential of ecologically complex environments.

## Metagenomic sequencing of gut microbiomes

As part of the MetaHIT (Metagenomics of the Human Intestine Tract) project, we collected faecal specimens from 124 healthy, overweight and obese individual human adults, as well as inflammatory bowel disease (IBD) patients, from Denmark and Spain (Supplementary Table 1). Total DNA was extracted from the faecal specimens and an average of 4.5 Gb (ranging between 2 and 7.3 Gb) of sequence was generated for each sample, allowing us to capture most of

*Bacteroides uniformis*  
*Alistipes putredinis*  
*Parabacteroides merdae*  
*Dorea longicatena*  
*Ruminococcus bromii* L2-63  
*Bacteroides caccae*  
*Clostridium* sp. SS2-1  
*Bacteroides thetaiotaomicron* VPI-5482  
*Eubacterium hallii*  
*Ruminococcus torques* L2-14  
Unknown sp. SS3 4  
*Ruminococcus* sp. SR1 5  
*Faecalibacterium prausnitzii* SL3 3  
*Ruminococcus lactaris*  
*Collinsella aerofaciens*  
*Dorea formicigenerans*

*Eubacterium rectale* M104 1  
*Bacteroides xylanisolvens* XB1A  
*Coprococcus comes* SL7 1  
*Bacteroides* sp. D1  
*Bacteroides* sp. D4  
*Eubacterium ventriosum*  
*Bacteroides dorei*  
*Ruminococcus obeum* A2-162  
*Subdoligranulum variabile*  
*Bacteroides capillosus*  
*Streptococcus thermophilus* LMD-9  
*Clostridium leptum*  
*Holdemanella filiformis*  
*Bacteroides stercoris*  
*Coprococcus eutectus*  
*Clostridium* sp. ME2 1  
*Bacteroides eggertii*  
*Butyrivibrio crossotus*  
*Bacteroides finegoldii*  
*Parabacteroides johnsonii*  
*Clostridium* sp. L2-50  
*Clostridium nexile*  
*Bacteroides pectinophilus*  
*Anaerotruncus colihominis*  
*Ruminococcus gnavus*  
*Bacteroides intestinalis*  
*Bacteroides fragilis* 3\_1\_12  
*Clostridium asparagiforme*  
*Enterococcus faecalis* TX0104  
*Clostridium scindens*  
*Blautia hansenii*



# Genome-Wide Mapping of in Vivo Protein-DNA Interactions

David S. Johnson,<sup>1\*</sup> Ali Mortazavi,<sup>2\*</sup> Richard M. Myers,<sup>1†</sup> Barbara Wold<sup>2,3†</sup>

In vivo protein-DNA interactions connect each transcription factor with its direct targets to form a gene network scaffold. To map these protein-DNA interactions comprehensively across entire mammalian genomes, we developed a large-scale chromatin immunoprecipitation assay (ChIPSeq) based on direct ultrahigh-throughput DNA sequencing. This sequence census method was then used to map in vivo binding of the neuron-restrictive silencer factor (NRSF; also known as REST, for repressor element-1 silencing transcription factor) to 1946 locations in the human genome. The data display sharp resolution of binding position [ $\pm 50$  base pairs (bp)], which facilitated our finding motifs and allowed us to identify data also have high sensitivity and specificity (area  $\geq 0.96$ ) and statistical confidence candidate interactions. These include pancreatic islet cell development.

Although much is known about transcription factor binding and action at genes, far less is known about the position and function of entire transcription interactomes, especially for organisms with large genomes. Now that human, mouse, and large genomes have been sequenced possible, in principle, to measure how transcription factor is deployed across the genome for a given cell type and physiological condition. Such measurements are important for systems-level studies because they provide a global map of candidate gene network input connections. These direct physical interactions between transcription factors or cofactors and the

<sup>1</sup>Department of Genetics, Stanford University School of Medicine, Stanford, CA, 94305-5120, USA. <sup>2</sup>Biology Division, California Institute of Technology, Pasadena, CA 91125, USA. <sup>3</sup>California Institute of Technology Beckman Institute, Pasadena, CA 91125, USA.

\*These authors contributed equally to this work.  
†To whom correspondence should be addressed. E-mail: woldb@its.caltech.edu (B.W.); myers@shgc.stanford.edu (R.M.M.)

## ChIP-Seq RIP-Seq

...

completeness, and high binding-site resolution. These data-quality and depth issues dictate whether primary gene network structure can be inferred with reasonable certainty and comprehensiveness, and how effectively the data can be used to discover binding-site motifs by computational methods. For these purposes, statistical robustness, sampling depth across the genome, absolute signal and signal-to-noise ratio must be good enough to detect nearly all in vivo binding locations for a regulator with minimal inclusion of false-positives. A further challenge in genomes large or small is to map factor-binding sites with high positional resolution. In addition to making com-

putational discovery of binding motifs feasible, this dictates the quality of regulatory site annotation relative to other gene anatomy landmarks, such as transcription start sites, enhancers, introns and exons, and conserved noncoding features (2). Finally, if high-quality protein-DNA interactome measurements can be performed routinely and at reasonable cost, it will open the way to detailed studies of interactome dynamics in response to specific signaling stimuli or genetic mutations. To address these issues, we turned to ultrahigh-throughput DNA sequencing to gain sampling power and applied size selection on immuno-enriched DNA to enhance positional resolution.

shown here differs from ChIP methods such as ChIP (1); ChIP-SAGE (4) in design, data design is simple (Fig. ChIP-Pet, it involves no array features. Unlike microarray of single-copy sites in ChIP-Seq assay (5), ChIP-Seq would be array features. Unlike microarray of single-copy sites in ChIP-Seq assay (5), ChIP-Seq would be array features. Unlike microarray of single-copy sites in ChIP-Seq assay (5), ChIP-Seq would be array features.

and so avoids constraints imposed by array hybridization chemistry, such as base composition constraints related to  $T_m$ , the temperature at which 50% of double-stranded DNA or DNA-RNA hybrids is denatured; cross-hybridization; and secondary structure interference. Finally, ChIP-Seq is feasible for any sequenced genome, rather than being restricted to species for which whole-genome tiling arrays have been produced.

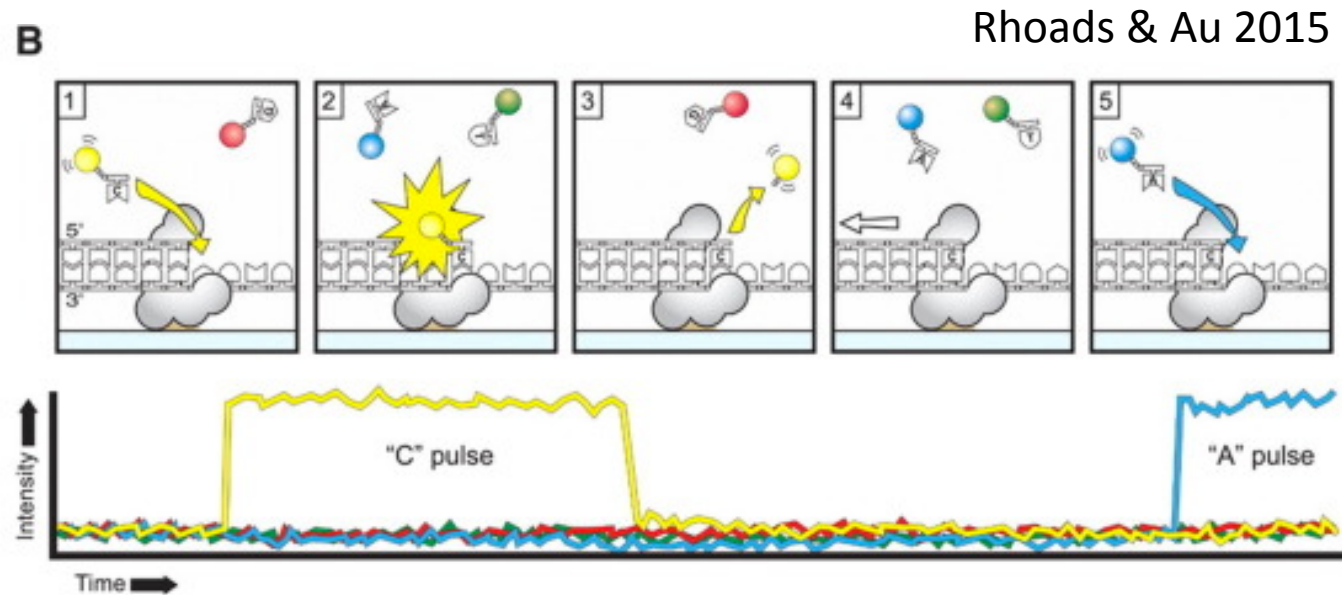
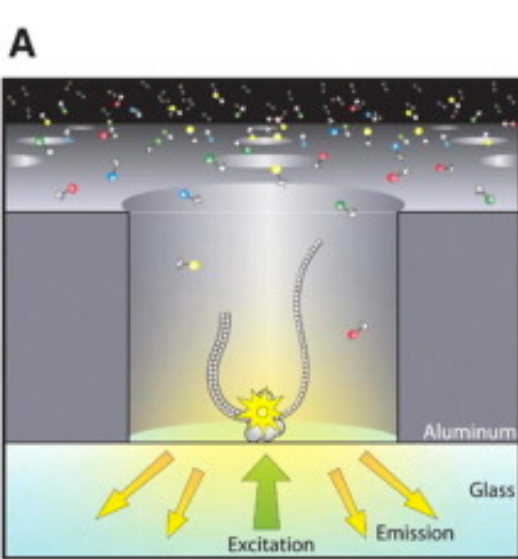
ChIP-Seq illustrates the power of new sequencing platforms, such as those from Solexa/Illumina and 454, to perform sequence census counting assays. The generic task in these applications is to identify and quantify the molecular



# Single molecule sequencing

- Sometimes called “third generation sequencing”
- Long read lengths (like Sanger sequencing) and high throughput (like second-generation sequencing)

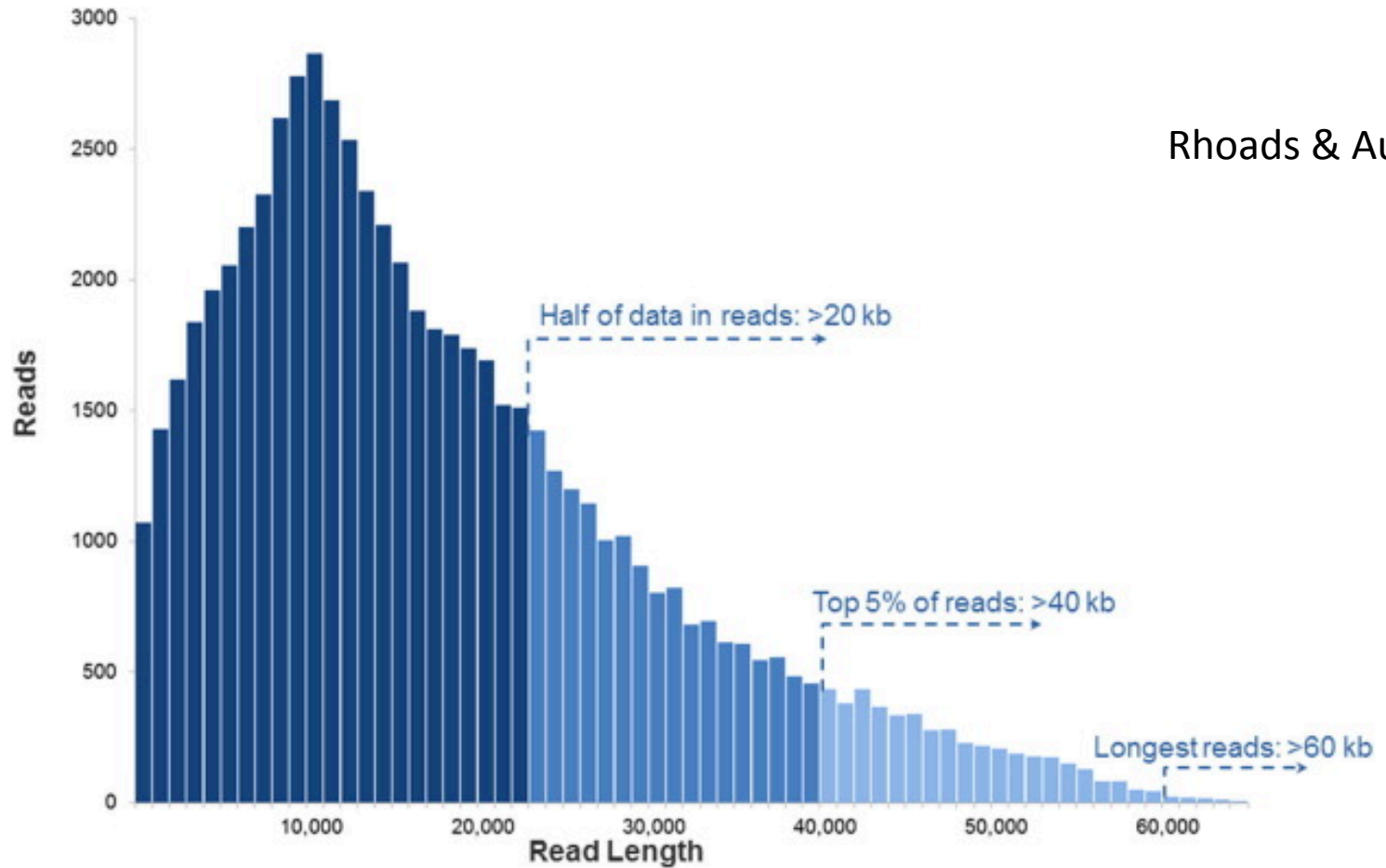
# ~2013: PacBio SMRT sequencing



Similar to 454 sequencing, but optics can detect the flash of a **single** dNTP being incorporated

# PacBio SMRT sequencing

Rhoads & Au 2015






Reads are **much** longer than Illumina reads



## RESEARCH ARTICLE

## An improved *Plasmodium cynomolgi* genome assembly reveals an unexpected methyltransferase gene expansion [version 1; referees: 2 approved]

Erica M Pasini<sup>1</sup>, Ulrike Böhme<sup>2</sup>, Gavin G. Rutledge<sup>1</sup> <sup>2</sup>,  
Annemarie Voorberg-Van der Wel<sup>1</sup>, Mandy Sanders<sup>2</sup>, Matt Berriman<sup>1</sup> <sup>2</sup>,  
Clemens HM Kocken<sup>1</sup>, Thomas D. Otto<sup>1</sup> <sup>2</sup>

<sup>1</sup>Biomedical Primate Research Centre, Rijswijk, Lange Kleiweg 161, 2288GJ Rijswijk, Netherlands

<sup>2</sup>The Wellcome Trust Sanger Institute, Hinxton, Cambridge, CB10 1SA, UK

**v1** First published: 16 Jun 2017, 2:42 (doi: 10.12688/wellcomeopenres.11864.1)  
Latest published: 16 Jun 2017, 2:42 (doi: 10.12688/wellcomeopenres.11864.1)

### Abstract

**Background:** *Plasmodium cynomolgi*, a non-human primate malaria parasite species, has been an important model parasite since its discovery in 1907. Similarities in the biology of *P. cynomolgi* to the closely related, but less tractable, human malaria parasite *P. vivax* make it the model parasite of choice for liver biology and vaccine studies pertinent to *P. vivax* malaria. Molecular and genome-scale studies of *P. cynomolgi* have relied on the current reference genome sequence, which remains highly fragmented with 1,649 unassigned scaffolds and little representation of the subtelomeres.

**Methods:** Using long-read sequence data (Pacific Biosciences SMRT technology), we assembled and annotated a new reference genome sequence, PcyM, sourced from an Indian rhesus monkey. We compare the newly assembled genome sequence with those of several other *Plasmodium* species, including a re-annotated *P. coatneyi* assembly.



**Results:** The new PcyM genome assembly is of significantly higher quality than the existing reference, comprising only 56 pieces, no gaps and an improved average gene length. Detailed manual curation has ensured a comprehensive annotation of the genome with 6,632 genes, nearly 1,000 more than previously attributed to *P. cynomolgi*. The new assembly also has an improved representation of the subtelomeric regions, which account for nearly 40% of the sequence. Within the subtelomeres, we identified more than 1300 *Plasmodium* interspersed repeat (*pir*) genes, as well as a striking expansion of 36 methyltransferase pseudogenes that originated from a single copy on chromosome 9.

**Conclusions:** The manually curated PcyM reference genome sequence is an important new resource for the malaria research community. The high quality and contiguity of the data have enabled the discovery of a novel expansion of methyltransferase in the subtelomeres, and illustrates the new comparative genomics capabilities that are being unlocked by complete reference genomes.

### Keywords

*P. cynomolgi*, PacBio assembly, *P. coatneyi*, methyltransferase

### Open Peer Review

Referee Status:  

Invited Referee  
1

version 1  
published  
16 Jun 2017  
 report

- 1 Aaron R. Jex, Walter and Eliza Institute of Medical Research, A
- 2 Richárd Bártfai, Radboud Univ Netherlands

### Discuss this article

Comments (0)

Better resolution of genomic sequence (especially in repeat/expansion regions)

# Defining a personal, allele-specific, and single-molecule long-read transcriptome

Hagen Tilgner<sup>a,1</sup>, Fabian Grubert<sup>a,1</sup>, Donald Sharon<sup>a,b,1</sup>, and Michael P. Snyder<sup>a,2</sup>

<sup>a</sup>Department of Genetics, Stanford University, Stanford, CA 94305-5120; and <sup>b</sup>Department of Molecular, Cellular, and Developmental Biology, Yale University, New Haven, CT 06511

Edited by Sherman M. Weissman, Yale University School of Medicine, New Haven, CT, and approved June 3, 2014 (received for review January 8, 2014)

Personal transcriptomes in which all of an individual's genetic variants (e.g., single nucleotide variants) and transcript isoforms (transcription start sites, splice sites, and polyA sites) are defined and quantified for full-length transcripts are expected to be important for understanding individual biology and disease, but have not been described previously. To obtain such transcriptomes, we sequenced the lymphoblastoid transcriptomes of three family members (GM12878 and the parents GM12891 and GM12892) by using a Pacific Biosciences long-read approach complemented with Illumina 101-bp sequencing and made the following observations. First, we found that reads representing all splice sites of a transcript are evident for most sufficiently expressed genes  $\leq 3$  kb and often for genes longer than that. Second, we added and quantified previously unidentified splicing isoforms to an existing annotation, thus creating the first personalized annotation to our knowledge. Third, we determined SNVs in a de novo manner and connected them to RNA haplotypes, including HLA haplotypes, thereby

for both parents of GM12878 (GM12891 and GM12892), we show that despite the higher error rate of the PacBio platform, single molecules can be attributed to the alleles from which they were transcribed, thereby generating accurate personal transcriptomes. This technique allows the assessment of biased allelic expression and isoform expression.

## Results

**Increased Full-Length Representation of RNA Molecules by Consensus Reads.** We sequenced  $\sim 711,000$  circular consensus reads (CCS) molecules from unamplified, polyA-selected RNA from the GM12878 cell line (see Fig. S1 for mapping strategy). We have recently shown that CCS often describe all splice sites of typical RNA molecules, although the success rate decreases as RNA length increases (11). The CCS we sequenced here were significantly longer (average 1,188 bp, maximum 6 kb) than

(Fig. 1A). Both comparisons showed a strong correlation between the number of molecules and the length of the molecules. The correlation was more pronounced for longer RNA molecules. The molecules were present in the HOP dataset. The number of molecules dropped significantly ( $P < 2.2 \times 10^{-16}$ ) in GM12891 and GM12892.

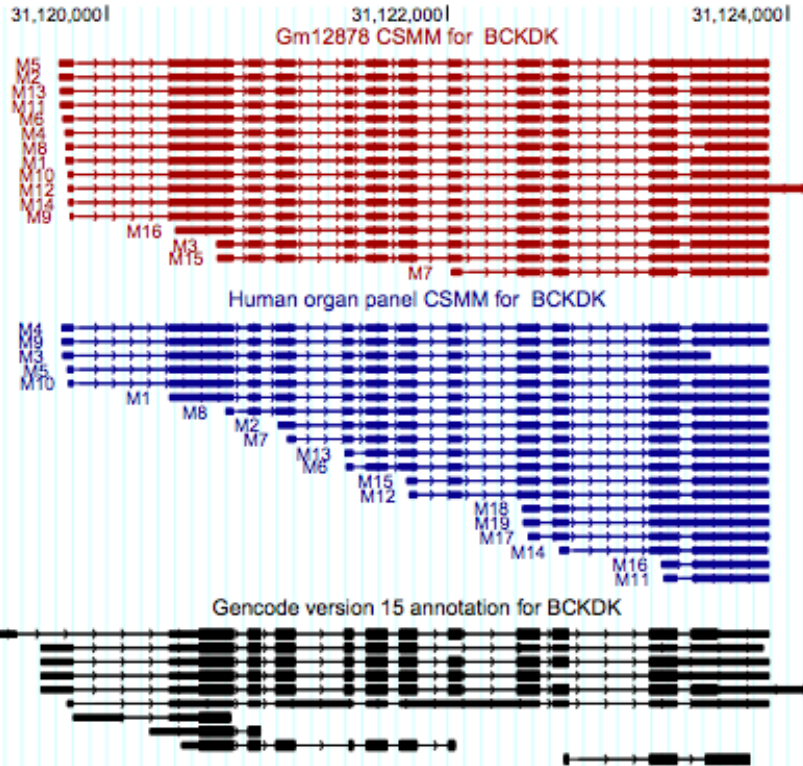
There can be thousands of molecules with two distinct variations (SNVs) in the same region. This is hard in RNA biology, although very useful for a long read. We show that one can find the allele it originates from structure variation, allele-specific splicing, and existing splicing existing. Identifying this enhanced paired-end reads.

Research: H.T., F.G., and D.S. wrote the paper. H.T. is a member of the advisory board of Personalis.

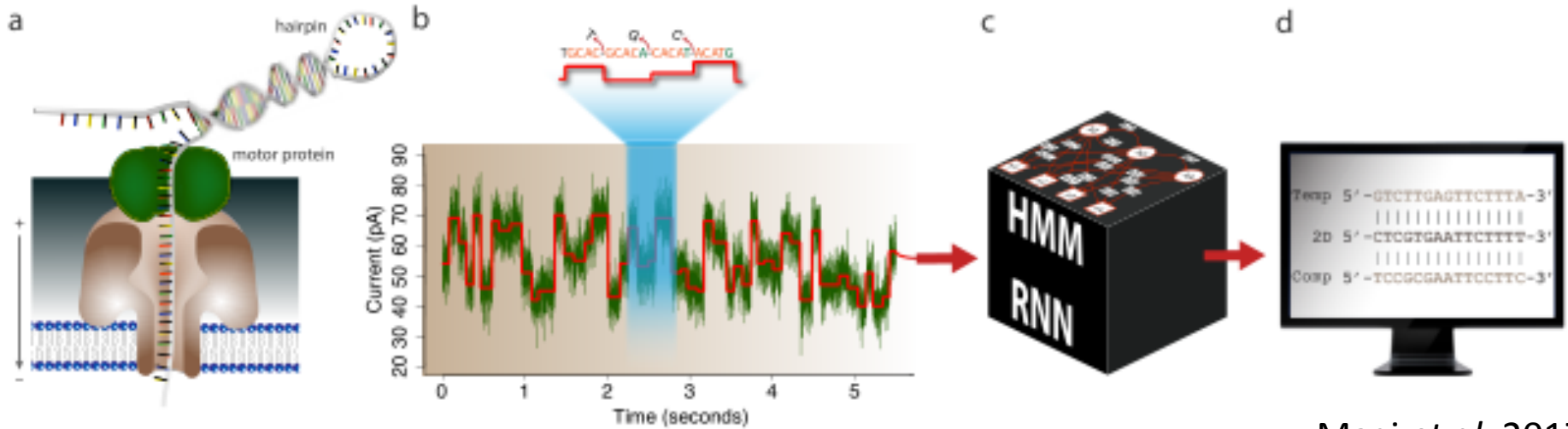
The data have been deposited in the NCBI SRA. The raw data are available at <http://www.ncbi.nlm.nih.gov/sra>.

Correspondence: [psnyder@stanford.edu](mailto:psnyder@stanford.edu).  
DOI: [10.1073/pnas.1317000110](https://doi.org/10.1073/pnas.1317000110).

Much better resolution of RNA sequences (e.g. splice variants)



# ~2014: Nanopore sequencing



Magi *et al.* 2017

**Not** sequencing by synthesis: direct reading of the sequence

# Nanopore sequencing



# Using long-read sequencing to detect imprinted DNA methylation

Scott Gigante<sup>1,a,c,✉</sup>, Quentin Gouil<sup>1,a,b</sup>, Alexis Lucattini<sup>c</sup>, Andrew Keniry<sup>a</sup>, Tamara Beck<sup>a</sup>, Matthew Tinning<sup>c</sup>, Lavinia Gordon<sup>c</sup>, Chris Woodruff<sup>a</sup>, Terence P. Speed<sup>a,d</sup>, Marnie E. Blewitt<sup>2,a,b</sup>, and Matthew E. Ritchie<sup>2,a,b,d,✉</sup>

<sup>1</sup>Equal first author

<sup>2</sup>Equal last author

<sup>a</sup>The Walter and Eliza Hall Institute of Medical Research, 1G Royal Parade, Parkville, 3052 Australia

<sup>b</sup>Department of Medical Biology, The University of Melbourne, Parkville, 3010 Australia

<sup>c</sup>Australian Genome Research Facility, 305 Grattan Street, Melbourne, 3000 Australia

<sup>d</sup>School of Mathematics and Statistics, The University of Melbourne, Parkville, 3010 Australia

<sup>e</sup>Department of Genetics, Yale University, New Haven, 06520 USA

Systematic variation in the methylation of cytosines at CpG sites plays a critical role in early development of humans and other mammals. Of particular interest are regions of differential methylation between parental alleles, as these often dictate monoallelic gene expression, resulting in parent of origin specific control of the embryonic transcriptome and subsequent development, in a phenomenon known as genomic imprinting.

Using long-read nanopore sequencing we show that, with an average genomic coverage of approximately ten, it is possible to determine both the level of methylation of CpG sites and the haplotype from which each read arises. The long-read property is exploited to characterise, using novel methods, both methylation and haplotype for reads that have reduced basecalling precision compared to Sanger sequencing. We validate the analysis both through comparison of nanopore-derived methylation patterns with those from Reduced Representation Bisulfite Sequencing data and through comparison with previously reported data.

Our analysis successfully identifies known imprinting control regions as well as some novel differentially methylated regions which, due to their proximity to hitherto unknown monoallelically expressed genes, may represent new imprinting control regions.

Nanopore sequencing | Differential methylation | Haplotyping | Imprinting | Long-read sequencing

Correspondence: [scott.gigante@yale.edu](mailto:scott.gigante@yale.edu), [mritchie@wehi.edu.au](mailto:mritchie@wehi.edu.au)

## Introduction

Methylation of the 5th carbon of cytosines (5mC or simply mC) is an epigenetic modification essential for normal mammalian development. Methylation differences between alleles contribute to establishing allele-specific expression patterns. As obtaining genome-wide haplotyped methylomes with short reads remains challenging, we evaluated the ability of long read, nanopore-based sequencing to improve allele-specific methylation analyses.

We apply the technique to the study of genomic imprinting, where differential expression of the maternal and paternal alleles in the offspring is at least partially set by the differential methylation (1–5). Imprinting is proposed to arise from the diverging interests of the maternal and paternal genes (6). In accordance with its primordial role in allocation of resources from the mother to the offspring, the placenta, along with the

brain, is the organ where parental conflict results in the most pronounced imprinted expression (7–9). We thus conduct a survey of differential methylation and expression in murine embryonic placenta.

Recent studies have increased the number of genes identified as subject to imprinting in mouse to about 200 (10–15). The cause of the differential expression between paternal and maternal alleles is only known for a subset of these genes; maternal histone marks can play a role (14), and in other cases it involves the differential methylation of adjacent regions (5). The differential methylation patterns may be established in the gametes and persist through the epigenetic reprogramming occurring after fertilisation (16). These differentially methylated regions (DMRs) are called primary DMRs, or imprinting control regions (ICRs). Alternatively, differential methylation may arise during development, perhaps as a downstream effect of differential expression, in which case the regions are called somatic or secondary DMRs (17).

Apart from the parent of origin of the allele, genetic differences can also be associated with differential methylation. In this case, F1 hybrids of distinct mouse strains will display DMRs between the alleles according to the strain of origin (18), and not the parent. Genetically determined DMRs can have profound effects on phenotype, for instance in humans by altering the expression of mismatch repair genes important in cancer (19). Therefore, we also investigate the link between DNA methylation and expression for strain-biased genes.

Reconstructing haplotyped methylomes necessitates the simultaneous measurement of DNA methylation and single-nucleotide polymorphisms (SNPs) differentiating the alleles. This can be achieved by deep sequencing of bisulfite-converted DNA on the Illumina platforms, although the short reads combined with the reduced complexity of the bisulfite-treated DNA make the process inefficient, meaning many regions with low SNP density remain unresolved. Long reads provided by third generation sequencing technologies can overcome the requirement of a high SNP density, while several methods allow the assessment of base modifications on native DNA (thus also avoiding the reduction in complexity associated with bisulfite conversion). These methods include: analysis of polymerase kinetics for PacBio SMRT se-



# DNA sequencing types: summary

Method	Read Length	Single pass error rate (%)	Reads per run	Time per run	Cost per million bases (USD)
Sanger (ABI)	600-1000	0.001	96	0.5 – 3h	500
454	700	1	1e6	23 h	8.57
Illumina (HiSeq)	2 x 125	0.1	8e9 (paired)	7 – 60 h	0.03
PacBio (RS II)	1-1.5e4	13	3.5-7.5e4	0.5 – 4 h	0.40-0.80
Oxford Nanopore (MinION)	2-5e3	38	1.1-4.7e4	50 h	6.44-17.90

# Data processing

Base calling (A/C/T/G)



Quality Scoring



Adapter trimming



FASTA/FASTQ files



Assembly



???



Alignment/mapping  
to reference

Performed by  
sequencing  
team

Performed by  
researcher

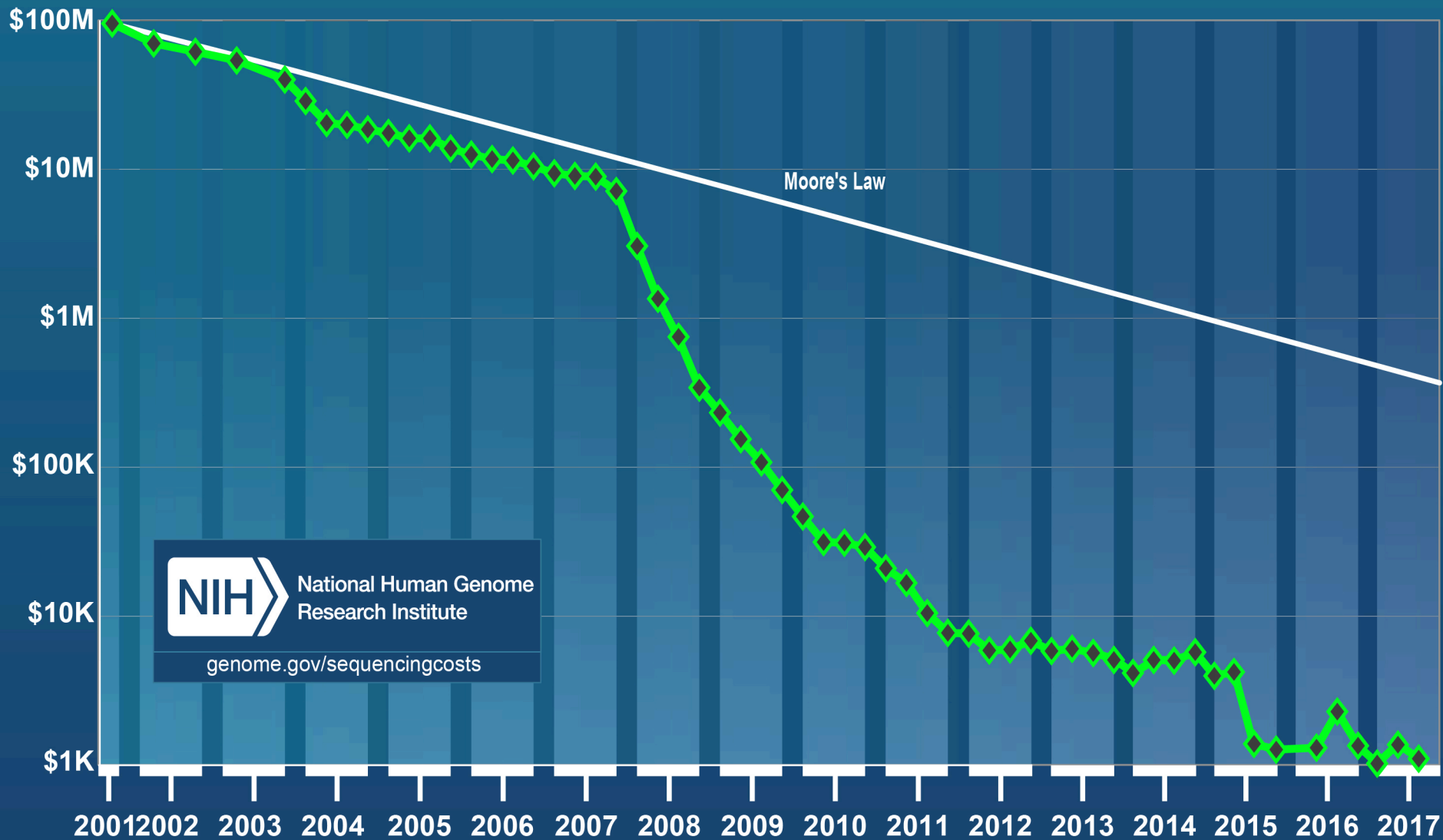




# Current challenges

- Repeats, especially homopolymers, are still difficult to sequence
  - Lower quality scores, less coverage
- Read lengths are still significantly shorter than some RNA transcripts
- Most technologies require amplification before sequencing, which can introduce errors and biases
- Most technologies don't provide information about base modifications
- Some technologies require large amounts of input material
- Cost, time, convenience

# Cost per Genome



# References

- Avery OT *et al.* 1944. Studies on the chemical nature of the substance inducing transformation of pneumococcal types: Induction of transformation desoxyribonucleic acid fraction isolated from pneumococcus type III. *Journal of Experimental Medicine* 79, 137–157
- Dahm R 2005. Freidrich Miescher and the discovery of DNA. *Developmental Biology* 278(2):274-288
- Griffith F. 1928. The significance of pneumococcal types. *Journal of Hygiene* 27: 113-159
- Johnson DS *et al.* Genome-wide mapping of in vivo protein-DNA Interactions. *Science* 316(5830): 1497-1502
- MacDonald ME *et al.* 1993. A novel gene containing a trinucleotide repeat that is expanded and unstable on Huntington's disease chromosomes. *Cell* 72(6): 971-983.
- Magi A *et al.* 2017. Nanopore sequencing data analysis: state of the art, applications and challenges. *Briefings in Bioinformatics*.
- Margulies M *et al.* 2005. Genome sequencing in microfabricated high-density picolitre reactors *Nature* 437: 376-380
- Pasini EA *et al.* 2017. An improved Plasmodium cynomolgi genome assembly reveals an unexpected methyltransferase gene expansion. *Wellcome Open Research* 2:42.
- Quick J *et al.* 2016. Real-time, portable genome sequencing for Ebola surveillance. *Nature* 530: 228-232
- Ratner L 1985. Complete nucleotide sequence of the AIDS virus, HTLV-III. *Nature* 313: 277-284.
- Rhoads A and Au KF. 2015. PacBio sequencing and its applications. *Genomics, Proteomics and Bioinformatics* 13(5):278-289
- Riordan JR *et al.* 1989. Identification of the cystic fibrosis gene: cloning and characterization of complementary DNA. *Science* 245(4922): 1066-1073.
- Sanger F *et al.* 1977. DNA sequencing with chain-terminating inhibitors. *PNAS* 74(12): 5463-5467
- Sanger F *et al.* 1977. Nucleotide sequence of bacteriophage  $\phi$ X174. *Nature* 265: 687-695
- Tilgner H *et al.* 2014. Defining a personal, allele-specific, and single-molecule long-read transcriptome. *PNAS* 111(27): 9869-9874.
- Valln HKM and Caldas C. 2011. The breast cancer genome – a key for better oncology. *BMC Cancer* 11: 501.
- Watson JD and Crick FHC. 1953. A structure for deoxyribose nucleic acid. *Nature* 171: 737–738