

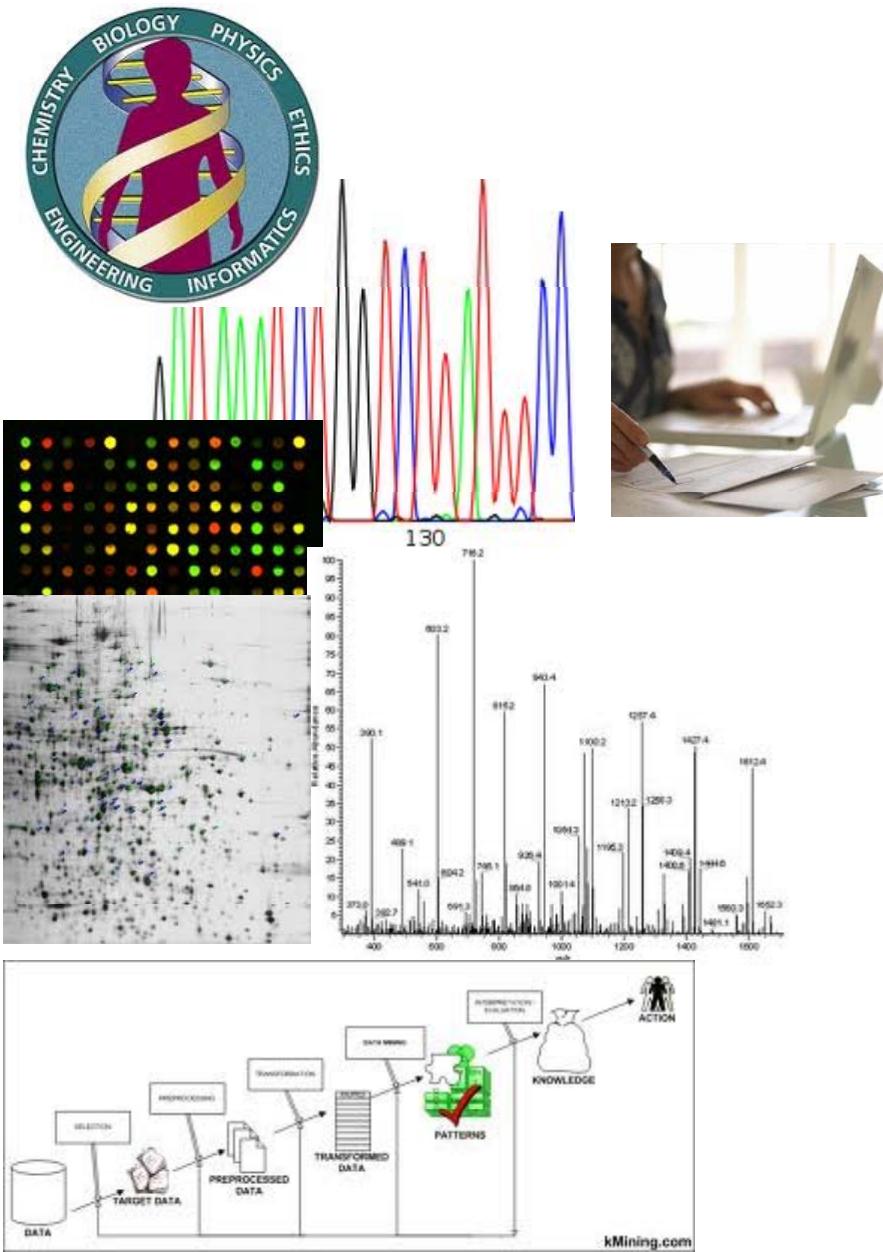
Curso de Posgrado: Alimentos Funcionales

BIONFORMATICA EN EL AMBITO DE LA NUTRIGENOMICA

Norma Paniego

Instituto de Biotecnología, CICVyA, CNIA INTA Castelar





CONCEPTOS GENERALES

NUTRIGENOMICA

NUTRIGENETICA

GENOMICA

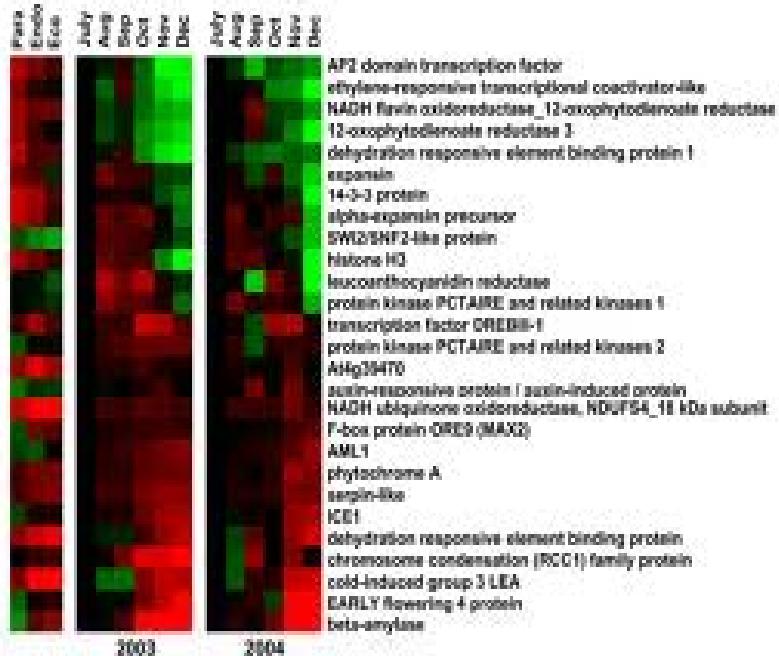
POST-GENOMICA

BASES DE DATOS

EXTRACCIÓN DEL CONOCIMIENTO

Nutrigenómica & Nutrigenética

La **nutrigenómica** es la ciencia que estudia la **expresión de los genes** en relación con la nutrición y el desarrollo de enfermedades asociadas a dicha expresión.



La **nutrigenética** es la ciencia que estudia como el **genotipo** del individuo interacciona con ingesta de nutrientes.



Proyectos genómicos

A horizontal timeline chart showing the progression of various genomic projects from 1995 to 2003. The timeline is represented by a blue background with vertical lines for each year. Key milestones are marked with icons and text. A red box highlights the final entry for 2003: "Finished version of human genome sequence completed".

- 1995:** HGP's human physical mapping goal achieved; First bacterial genome (*E. coli*) sequenced.
- 1996:** First human gene map established; Pilot projects for human genome sequencing begin in U.S.; first orchid genome sequenced; Yeast (*S. cerevisiae*) genome sequenced; U.S. Equal Employment Opportunity Commission issues policy on genetic discrimination in the workplace; HGP's mouse genetic mapping goal achieved; Bermuda principles for rapid and open data release established.
- 1997:** DOE forms Joint Genome Institute (JGI); ICHGR becomes NHGRI; *E. coli* genome sequenced; Genoscope (French National Genome Sequencing Center) founded.
- 1998:** Incorporation of 30,000 genes into human genome map; New five-year plan for the HGP in the U.S. published; RIKEN Genomic Sciences Center (Japan) established; Roundworm (*C. elegans*) genome sequenced; SNP initiative begins; GTGCT GTCCT sequence logo; Chinese National Human Genome Centers (in Beijing and Shanghai) established.
- 1999:** Full-scale human sequencing begins; Fruit fly (*D. melanogaster*) genome sequenced; Sequence of first human chromosome (chromosome 22) completed; Mustard cress (*A. thaliana*) genome sequenced.
- 2000:** Draft version of human genome sequence completed; President Clinton and Prime Minister Blair support free access to genome information; Executive order bans genetic discrimination in U.S. federal workplace.
- 2001:** Draft version of human genome sequence published; President Clinton and Prime Minister Blair support free access to genome information; *nature* magazine cover featuring the human genome.
- 2002:** Draft version of mouse genome sequence completed and published; Draft version of rat genome sequence completed; 10,000 full-length human cDNAs sequenced; Mammalian Gene Collection; Draft version of rice genome sequence completed and published.
- 2003:** Finished version of human genome sequence completed; HGP ends with all goals achieved; *nature* magazine cover featuring the finished human genome; "to be continued" note.

NCBI Entrez Genome Project

connection information discovery

Entrez Genome Gene Nucleotide Protein PopSet PubMed Taxonomy

Search: Genome Project Go Clear

Properties of Eukaryotic Genome Sequencing Projects

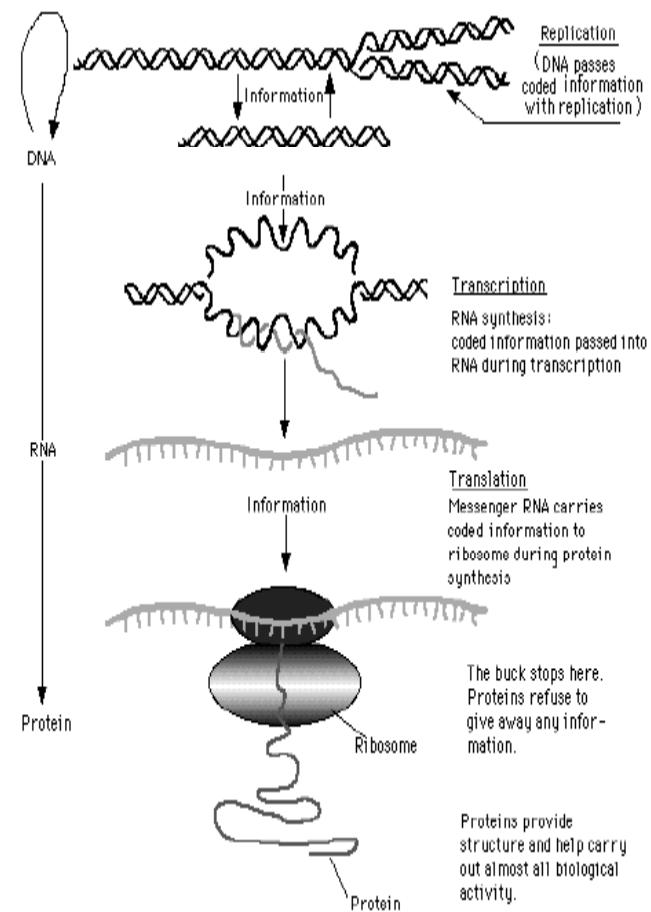
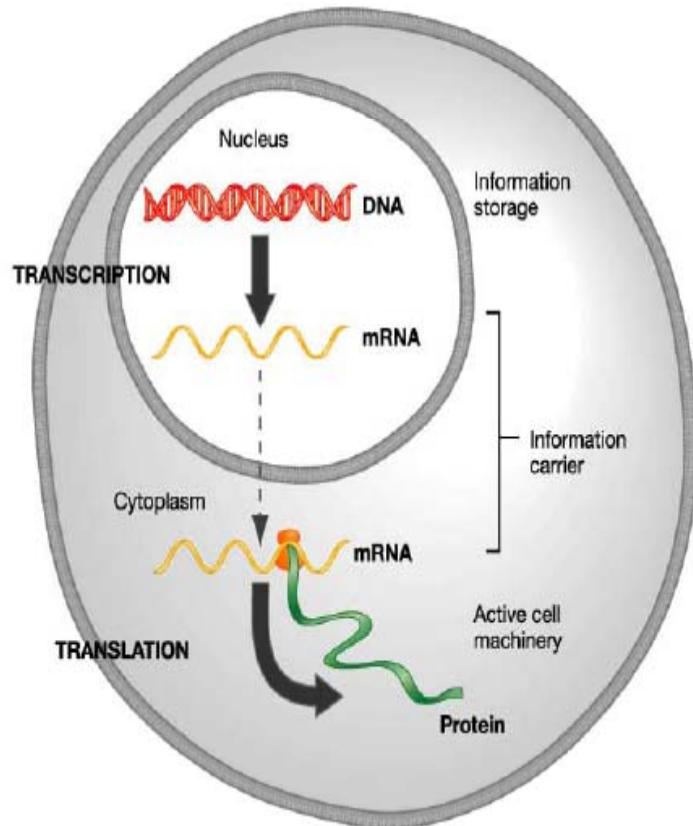
Organism Group: All Eukaryotes Sequencing Status: All Sequencing Method: All Go Reset

Abbreviations: GB - GenBank Accessions; PM - PubMed; R - RefSeq Accessions; G - Entrez Gene; T - Trace Archive; B - BLAST; M - Map Viewer; F - FTP Sites

482 Eukaryotic Genome Sequencing Projects Selected: Complete - 23, Assembly - 231, In Progress - 228 save

731 Eukaryotic Genome Sequencing Projects Selected: Complete - 23, Assembly - 327, In Progress - 381 save

El Dogma Central de la Biología Molecular



Niveles básicos de Información biológica

- **Genoma:** es la totalidad de la información genética que posee un organismo en particular.
- **Transcriptoma:** la parte del genoma que se expresa en una célula en una etapa específica de su desarrollo.
- **Proteoma:** la totalidad de las proteínas codificadas por un genoma.
- **Metaboloma:** la totalidad de pequeñas moléculas o metabolitos que se pueden encontrar en una muestra biológica, tal como un organismo.

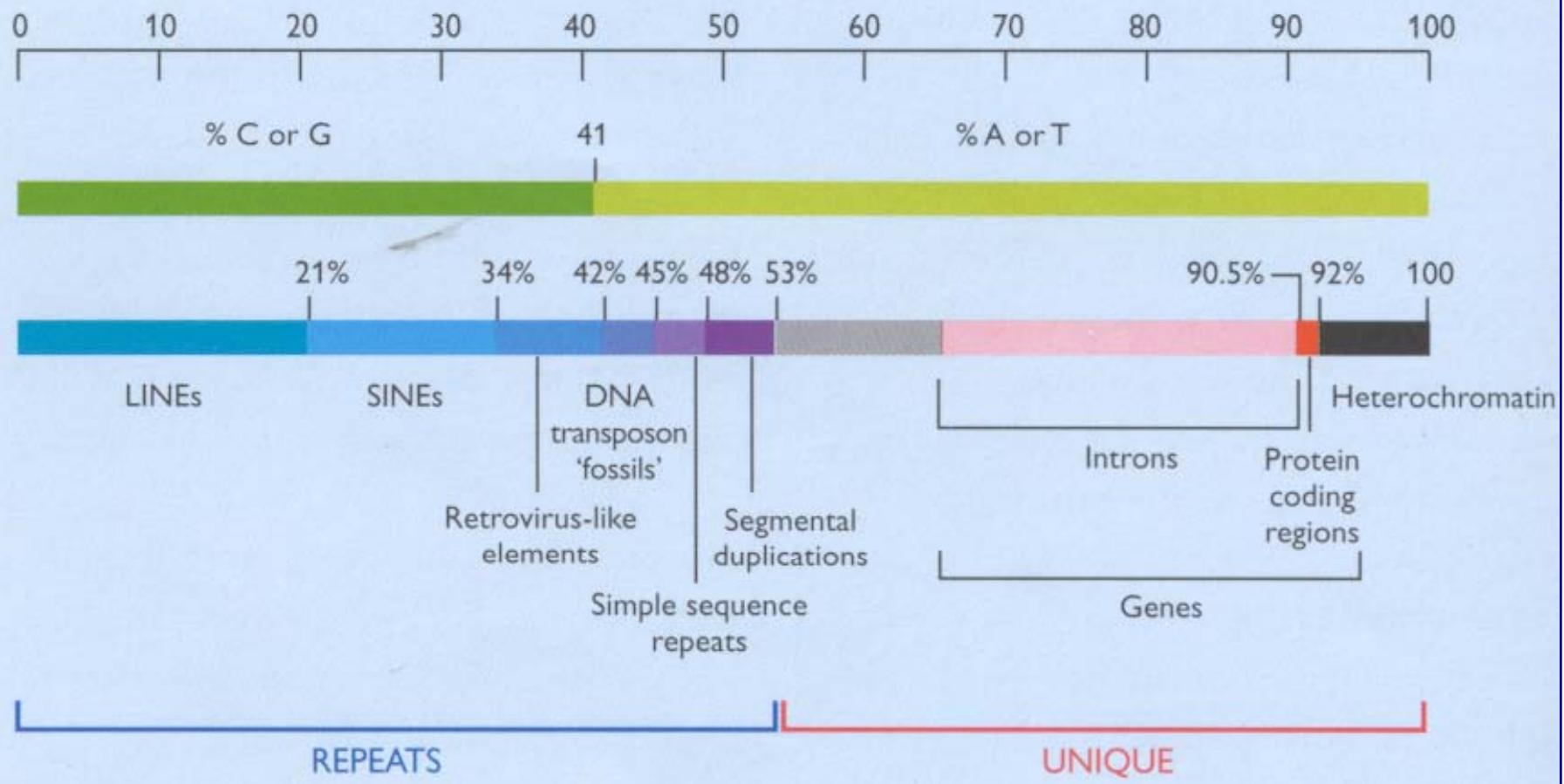
La era de la genómica

La genómica se ha desarrollado como consecuencia de los avances en Biología Molecular e Informática.

La introducción y popularización de las tecnologías de alta procesividad ha cambiado drásticamente la manera en que se abordan los problemas biológicos y se prueban las hipótesis.

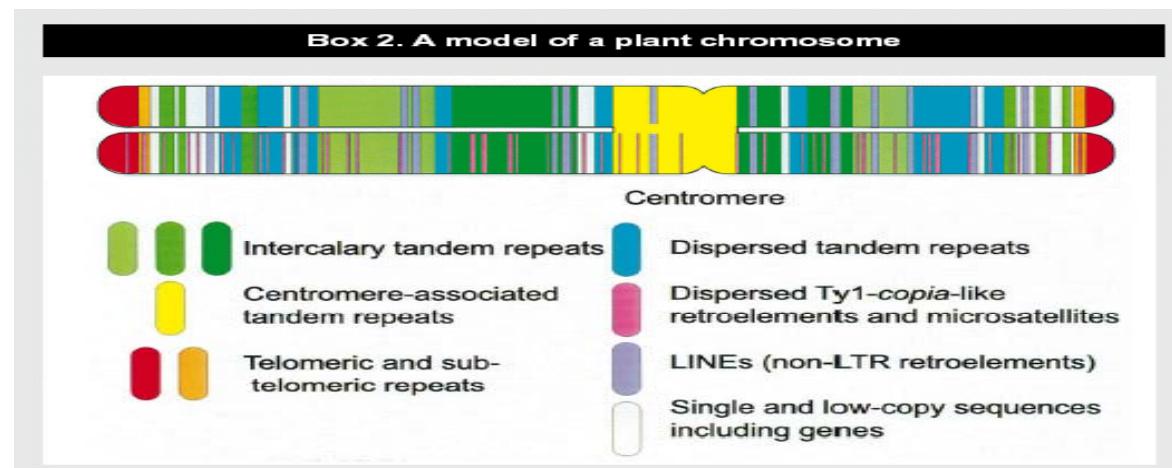
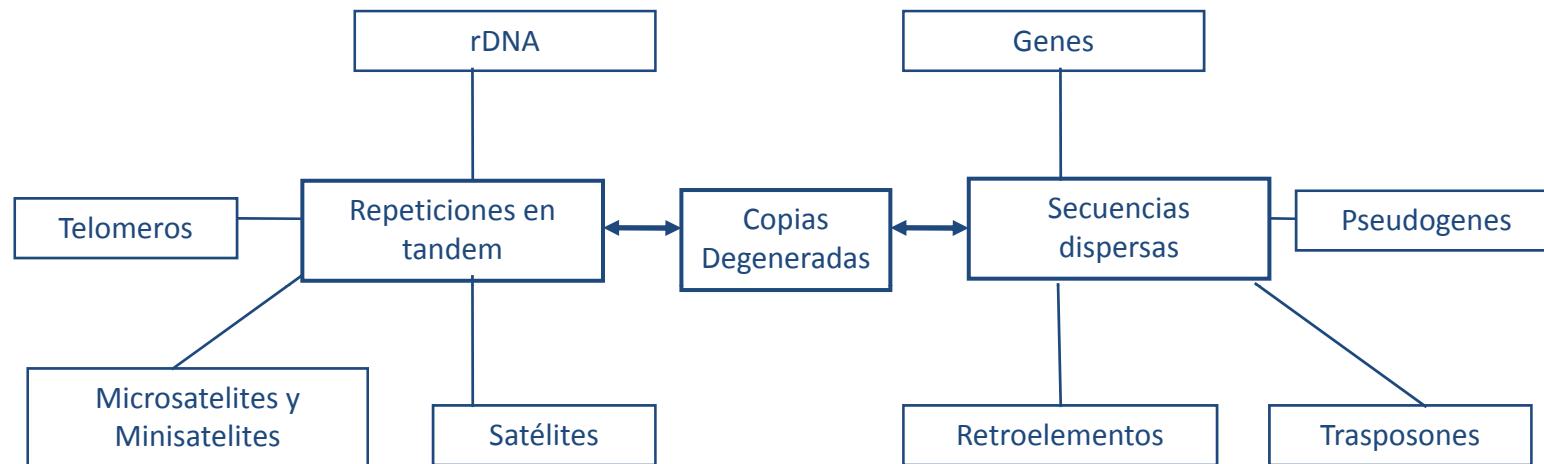
Inicialmente los proyectos genómicos se focalizaban en la secuencia completa del genoma, actualmente incluyen el análisis de la expresión y función de genes, proteínas, metabólitos entre otros.

Figure 1: The genome by numbers



[The Genome by Numbers](#), the Wellcome Trust.

Contenido de secuencias de ADN en el genoma vegetal

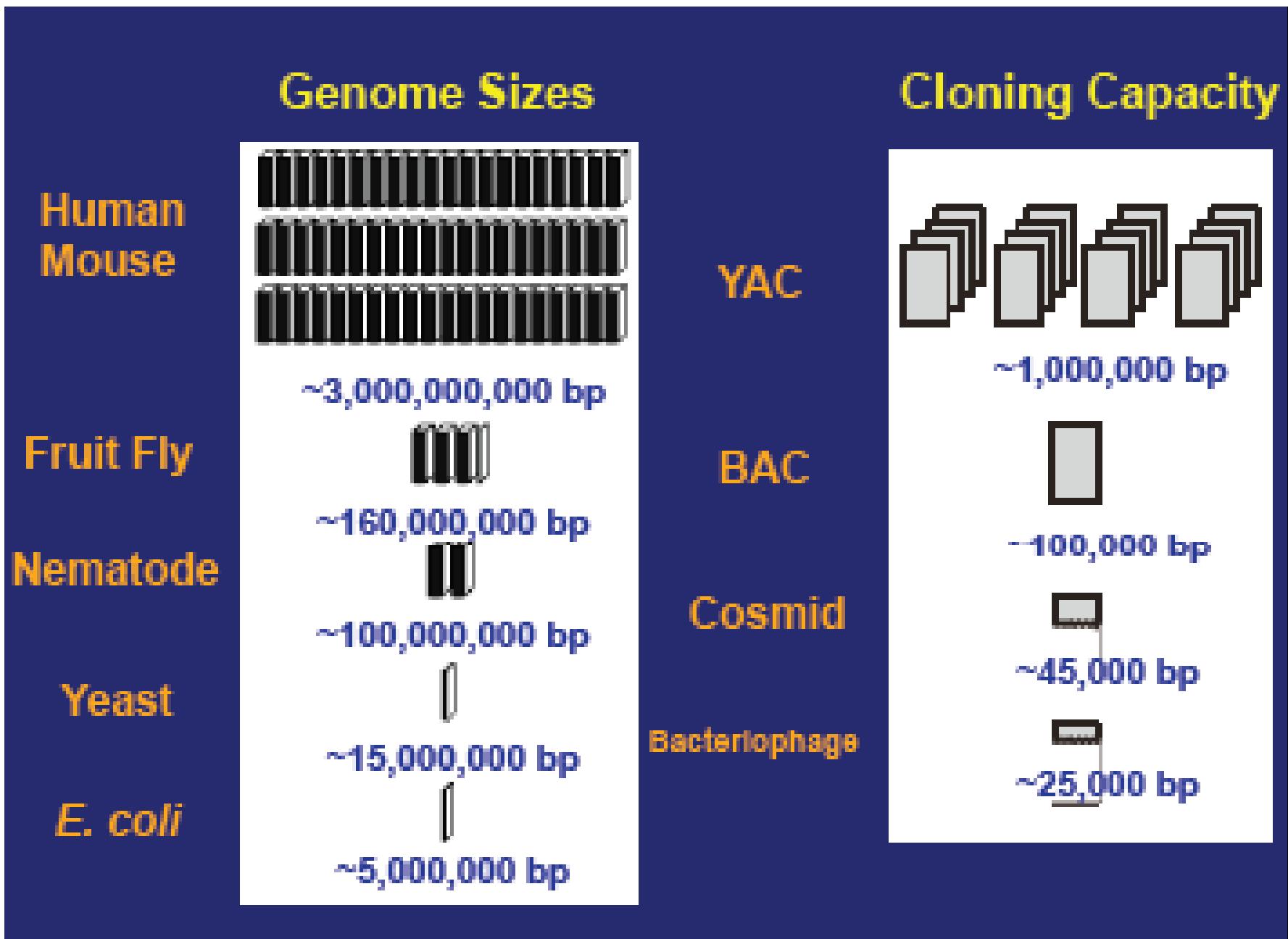


Schmid T y Heslop-Harrison JS (1998) Trends in Plant Sci., 3, 195-199

Single nucleotide variant	ATTGGCCTTAACCCCGATTATCAGGAT ATTGGCCTTAACCTCCGATTATCAGGAT
Insertion–deletion variant	ATTGGCCTTAACCCGATCCGATTATCAGGAT ATTGGCCTTAACCC---CCGATTATCAGGAT
Block substitution	ATTGGCCTTAACCCCCGATTATCAGGAT ATTGGCCTTAACAGTG GATTATCAGGAT
Inversion variant	ATTGGCCTTAAACCCCCGATTATCAGGAT ATTGGCCTTCCGGGGTTATTATCAGGAT
Copy number variant	ATTGGCCTTAGGCCTTAACCCCCGATTATCAGGAT ATTGGCCTTA-----ACCTCCGATTATCAGGAT

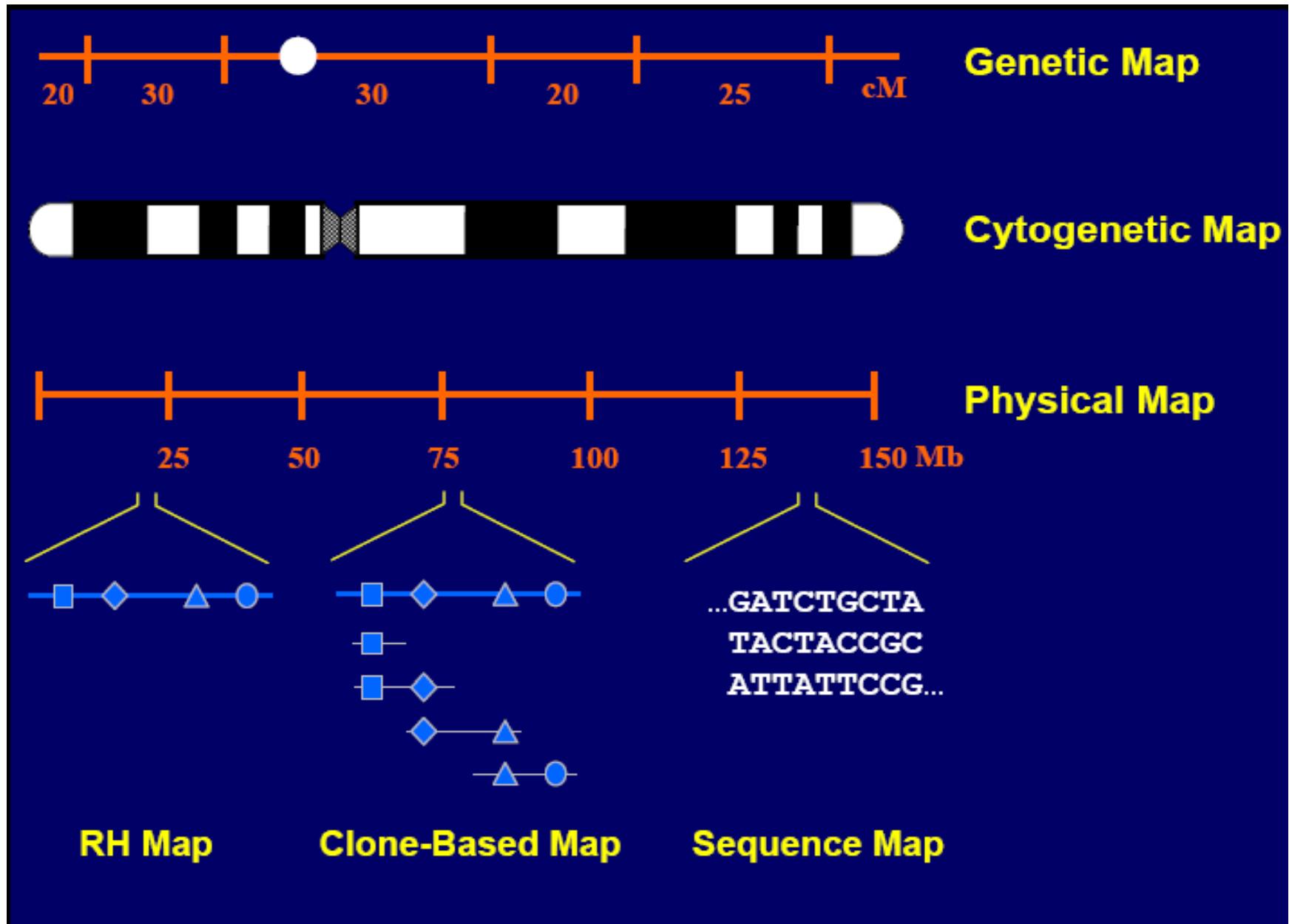
Structural variants

Figure 1 | Classes of human genetic variants. The nomenclature used to describe the various types of structural variants is not yet standard¹²¹. Here, the terminology used aims to describe the nucleotide composition of the variant and distinguish it from other types of variants. Single nucleotide variants are DNA sequence variations in which a single nucleotide (A, T, G or C) is altered. Insertion–deletion variants (indels) occur when one or more base pairs are present in some genomes but absent in others. They are generally composed of only a few bases but can be greater than 80 kb in length¹¹. Block substitutions describe cases in which a string of adjacent nucleotides varies between two genomes. An inversion variant is one in which the order of the base pairs is reversed in a defined section of a chromosome. A well-characterized inversion variant that has been described in humans involves a section of chromosome 17 in which a ~900 kb interval is in the reverse order in approximately 20% of individuals with Northern European ancestry¹²². Copy number variants occur when identical or nearly identical sequences are repeated in some chromosomes but not others. The largest copy number variant identified in the Venter genome¹¹ was almost 2 Mb in length.



The Genomic Landscape: *circa 2010, Eric Green*

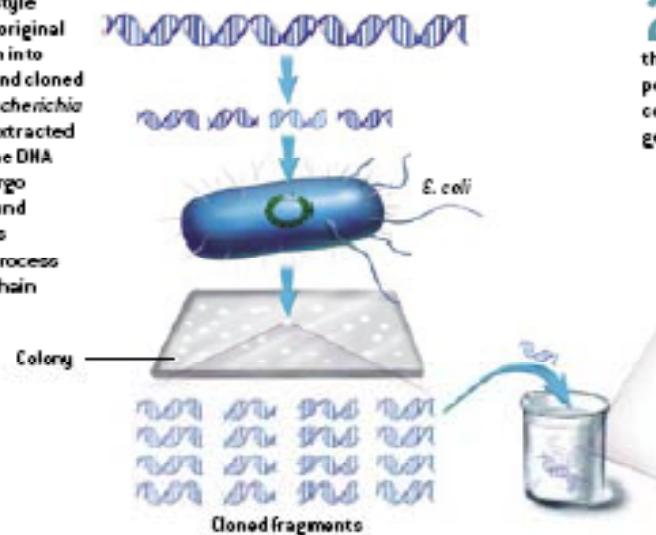
Como se estudian los genomas: técnicas de mapeo



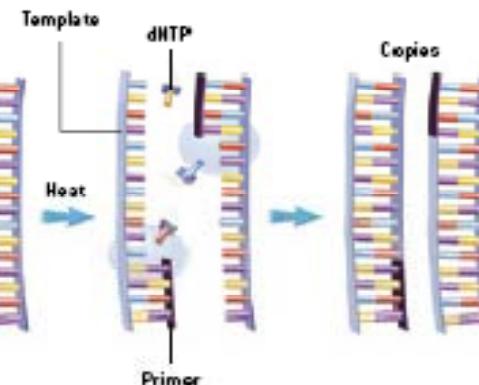
The Genomic Landscape: *circa 2010, Eric Green*

Método de secuenciación

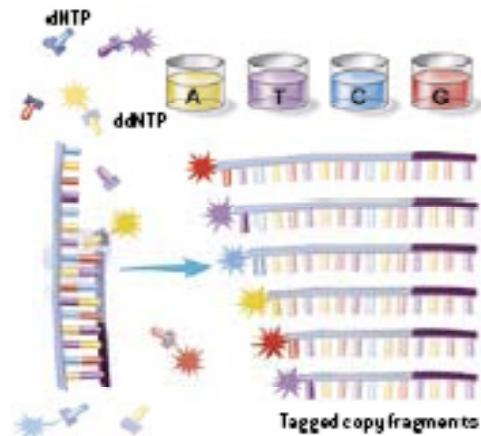
1 Before Sanger-style sequencing, an original DNA strand is broken into smaller fragments and cloned within colonies of *Escherichia coli* bacteria. Once extracted from the bacteria, the DNA fragments will undergo another massive round of copying, known as amplification, by a process called polymerase chain reaction (PCR).



2 During PCR, fragments are heated so they will separate into single strands. A short nucleotide sequence called a primer is then annealed to each original template. Starting at the primer, polymerase links free-floating nucleotides (called dNTPs) into new complementary strands. The process is repeated over and over to generate millions of copies of each fragment.



3 Single-stranded fragments are next tagged in a process similar to PCR but with fluorescently labeled terminators (ddNTPs) added to the mixture of primers, polymerase and dNTPs. Complementary strands are built until by chance a ddNTP is incorporated, halting synthesis. The resulting copy fragments have varying lengths and a tagged nucleotide at one end.



4 Capillary electrophoresis separates the fragments, which are negatively charged, drawing them toward a positively charged pole. Because the shortest fragments travel fastest, their order reflects their size and their ddNTP terminators can thus be read in template's base sequence. Laser light activates the fluorescent tags as the fragments pass through a 'Detection window', producing a color readout that is translated into a sequence.

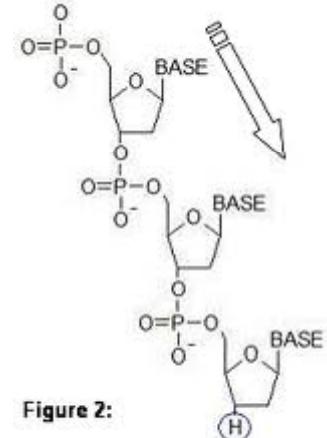
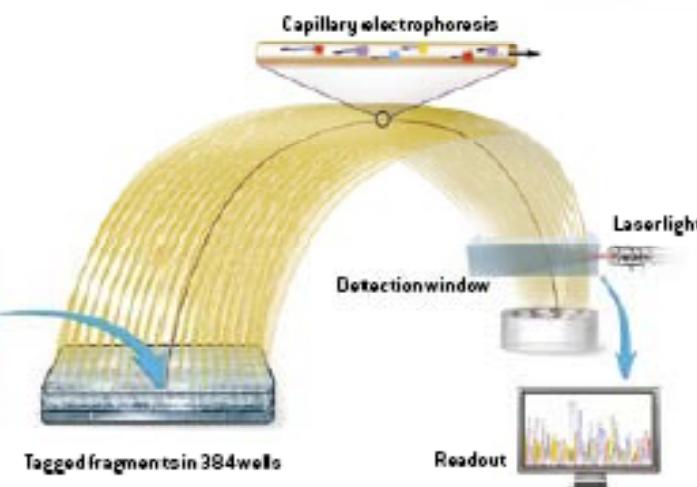


Figure 2:

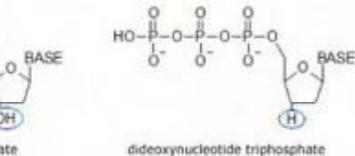
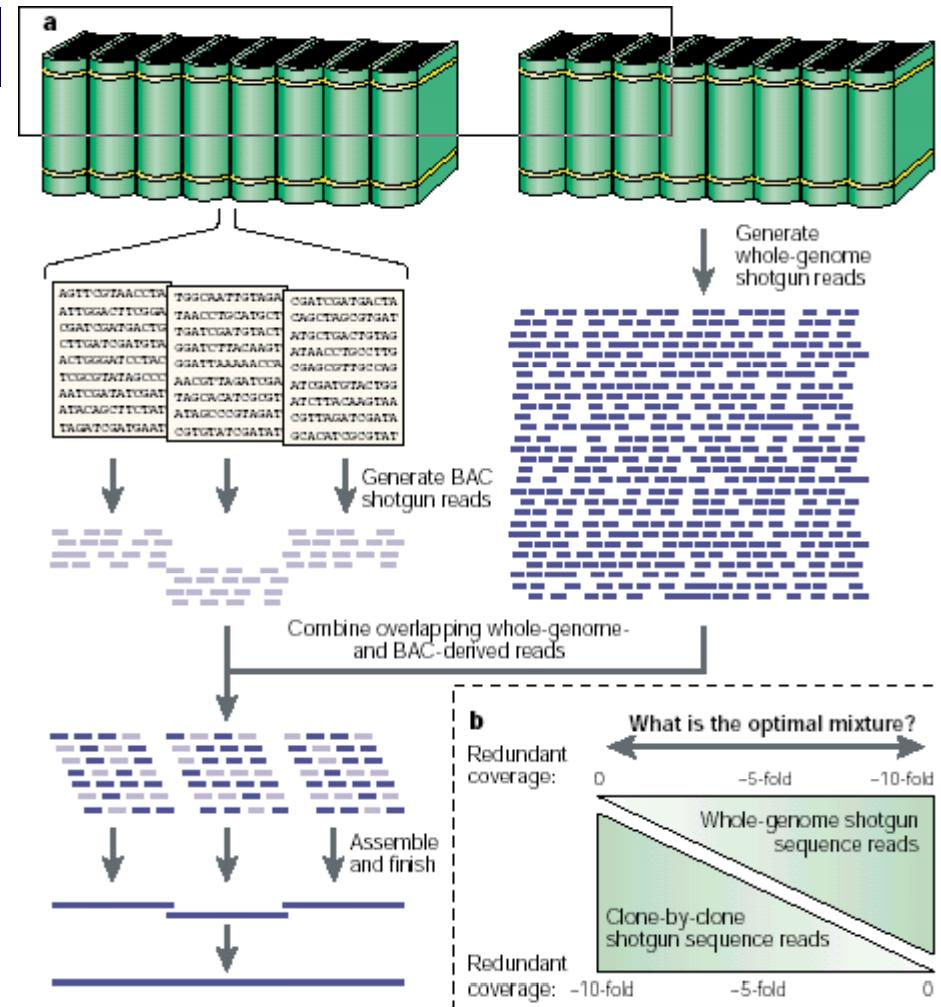
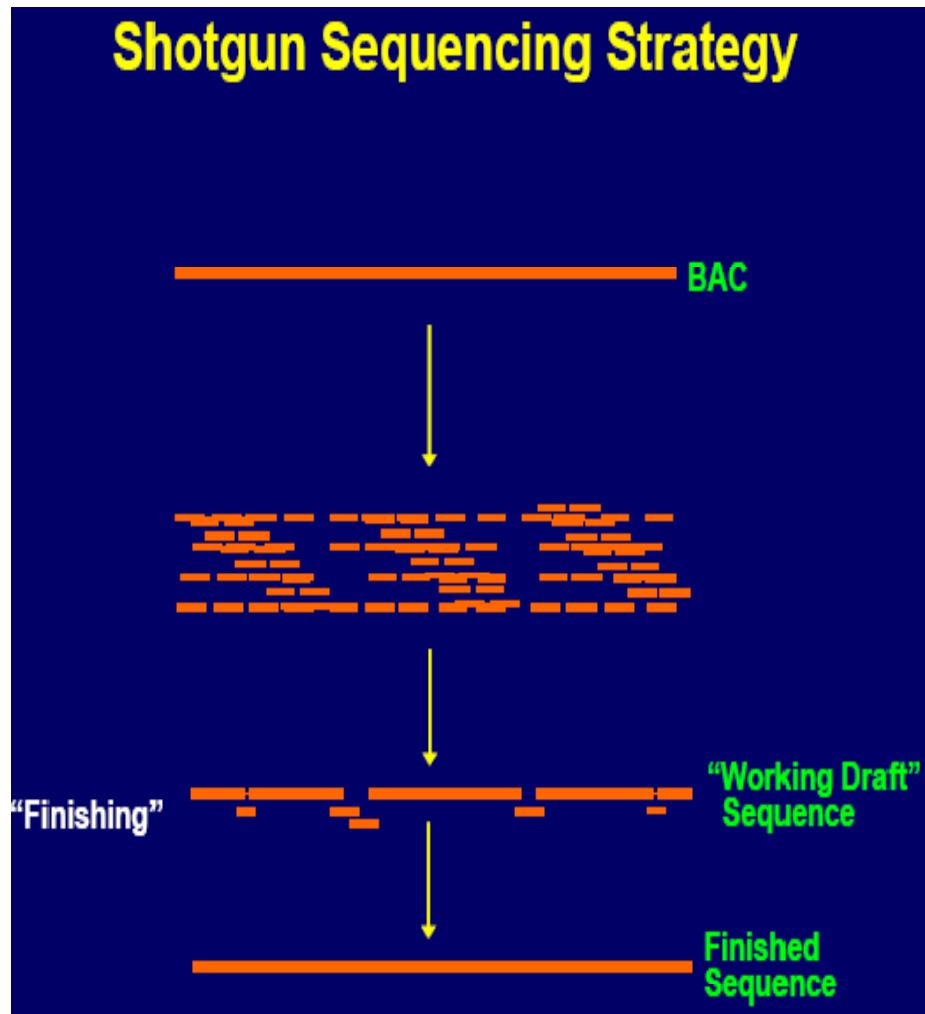


Figure 1:

Estrategias alternativas de secuenciación



The Genomic Landscape: *circa 2010, Eric Green*

Equipamiento automatizado

Human Genome Sequencing Centers



The Genomic Landscape: *circa 2010, Eric Green*

a) Sequence reads

Read 1 GACATACACATGG

Read 2 TCAATGGGGCTAA

Read 3 AGCAAGGACTTGTGACATACACATG

Read 4 ACACATGGAAATA

Read 5 GGGGTAATGATTGTCAC

Read 6 TGATTGTCACATA

Read 7 ATTGATGAGGCACCGA

Read 8 GTGACATACACATGATCAATGGGG

↓ Use computer to assemble sequence reads

↓

b) 7 ATTGATGAGGCACCGA

3 AGCAAGGACTTGTGACATACACATG

8 GTGACATACACATGATCAATGGGG

2 TCAATGGGGCTAA

5 GGGGTAATGATTGTCAC

6 TGATTGTCACATA

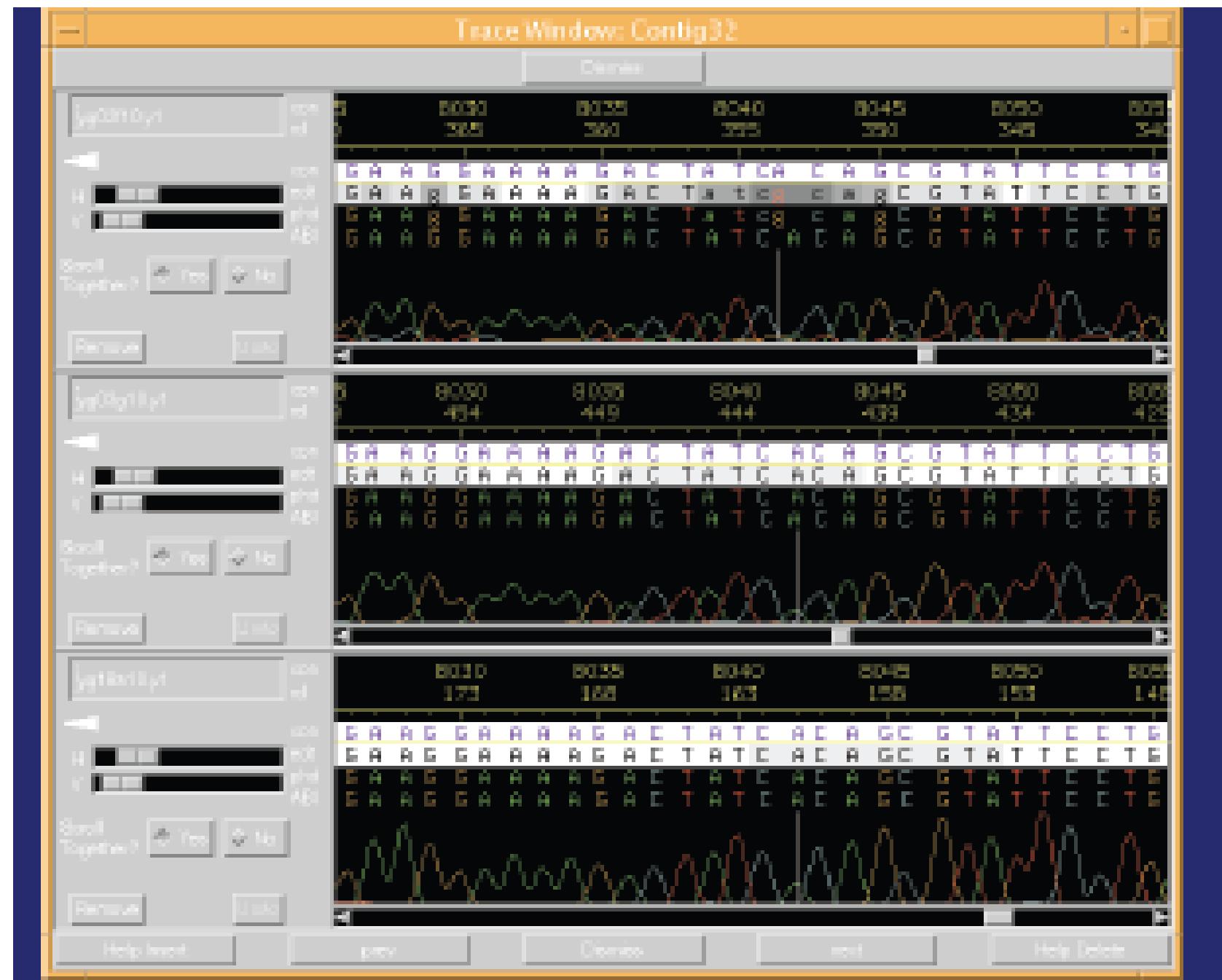
1 GACATACACATGG

4 ACACATGGAAATA

↓ Assembled sequence

↓

c) ATTGATGAGGCACCGACTTGTGACATACACATGATCAATGGGGCTAAATGATTGTCACATACACATGGAAATA



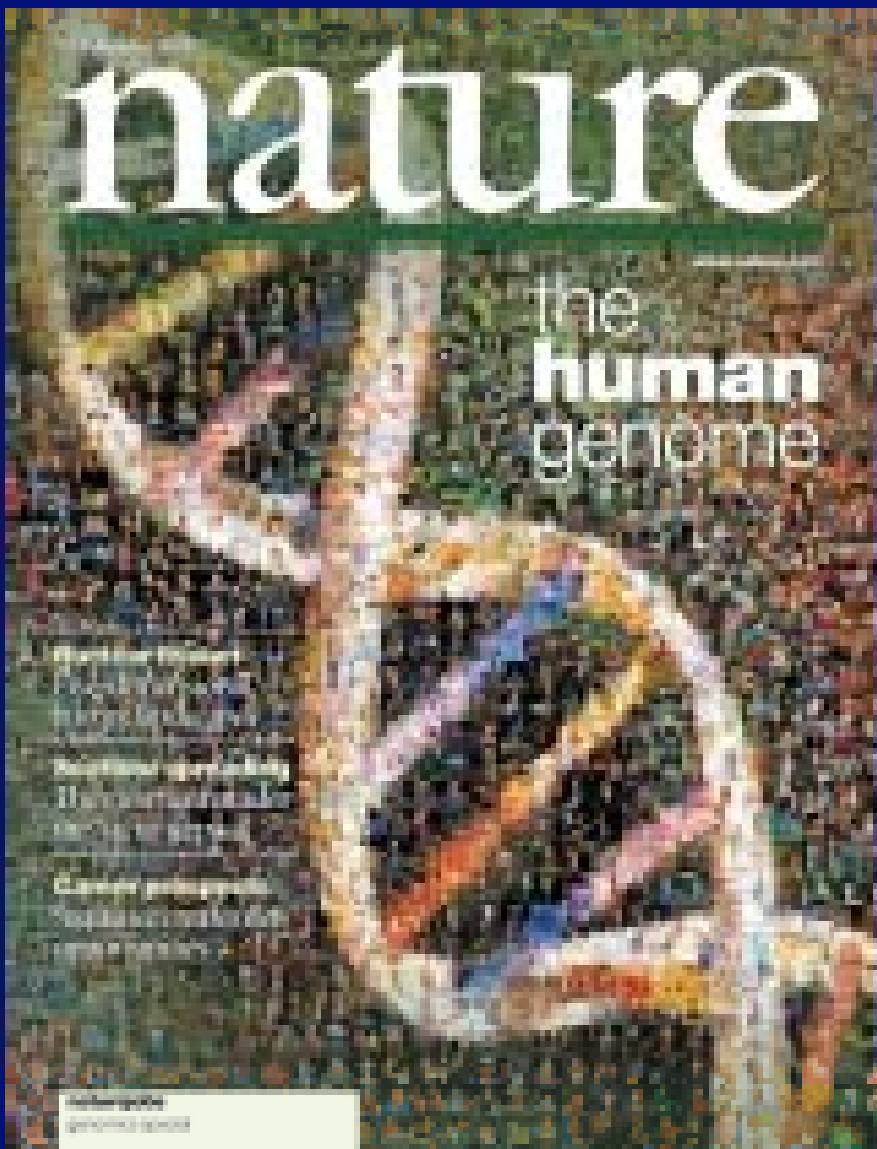
The Genomic Landscape: circa 2010, Eric Green

Sequence Finishing: Resolving Ambiguities

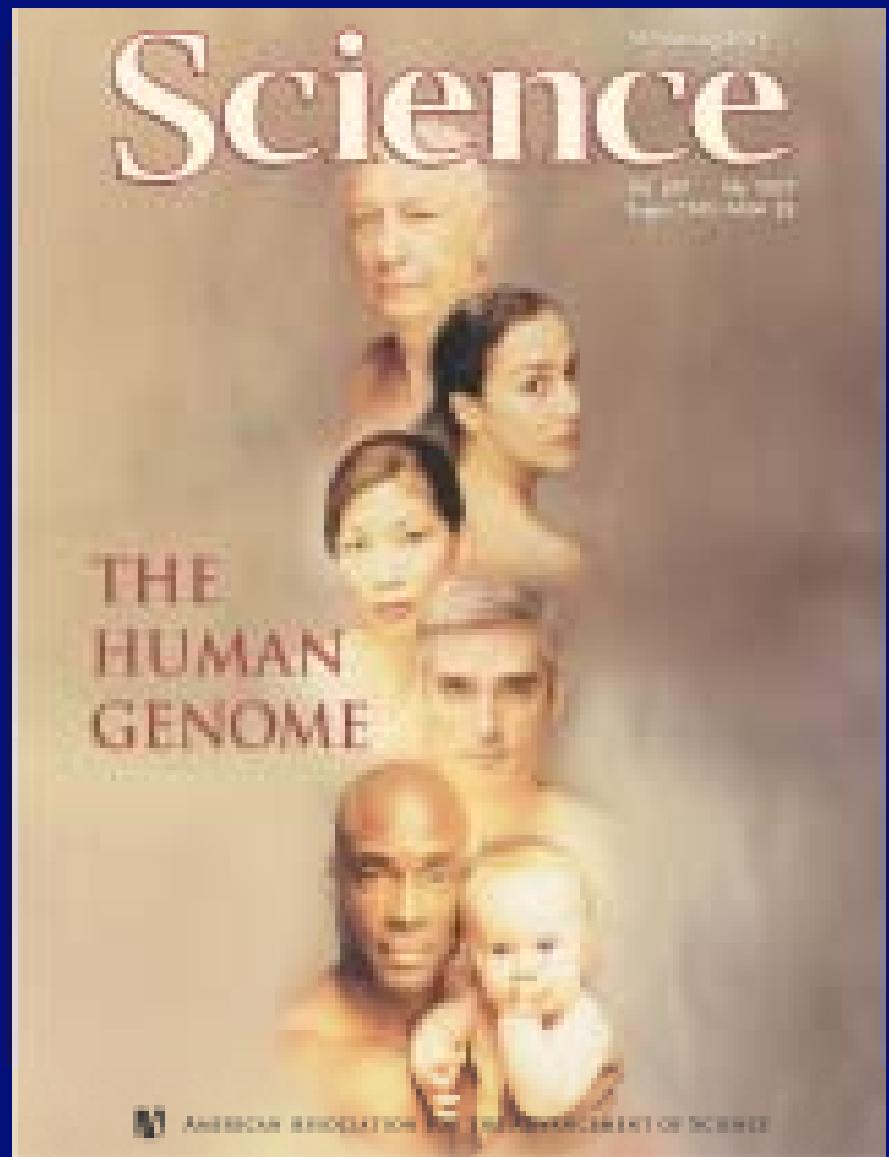


*** Sequence Finishing: Remains Relatively Expensive ***

February, 2001 Draft Sequence

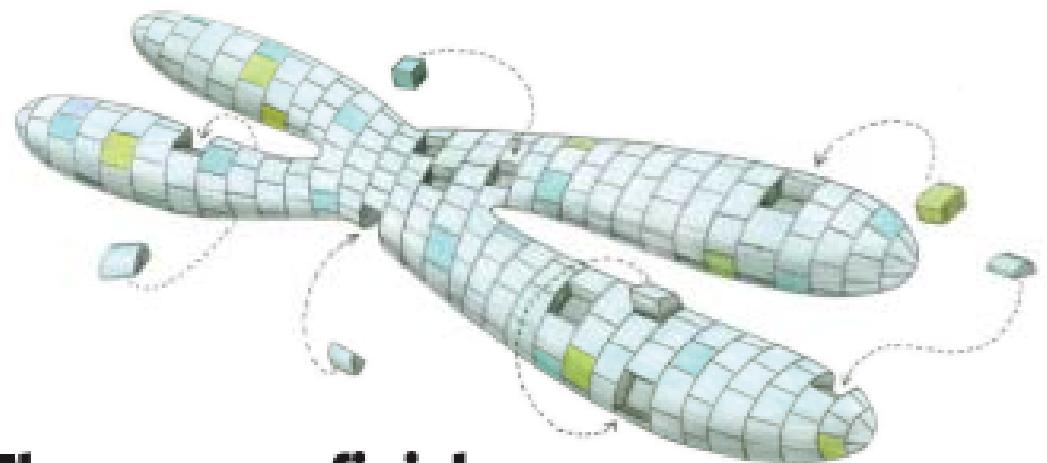
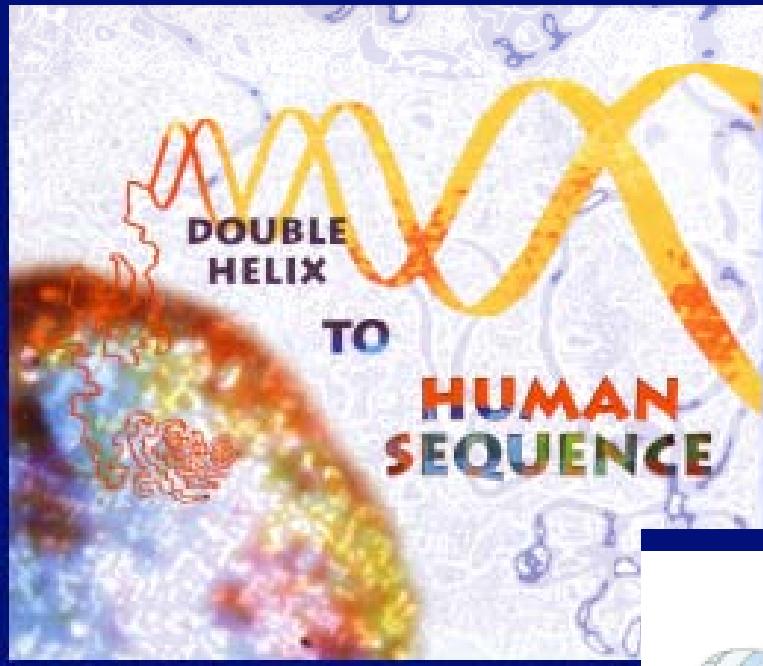


International Human Genome
Sequencing Consortium (2001)



Venter et al. (2001)
The Genomic Landscape: circa 2010, Eric Green

April, 2003 Completion



The genome finishers

Dedicated scientists are working hard to close the gaps, fix the errors and finally complete the human genome sequence. [Elie Delgin](#) looks at how close they are.

Nature (2009)

The Genomic Landscape: *circa 2010*, Eric Green

Mapping the Human Genome

~1990 to ~2000

The Human
Genome Project

Sequencing the Human Genome

~1998 to ~2003

Interpreting the Human Genome Sequence

~2003 to ???

Beyond
The Human
Genome Project

The Human Genome... by the Numbers

~5% of Human Genome Sequence is Constrained Across Mammals (and Presumed Functional)

5% of 3B Bases = ~150M Bases

Do NOT Yet Know the Position of these ~150M Functional Bases
Lower Bound for the Amount that is Functional

~1.5% Encodes for Protein (Genes)

Corresponds to ~18-22K Genes

Many More than ~22K Different Proteins

Good Inventory at Present

~3.5% Functional But Non-Coding

Gene Regulatory Elements

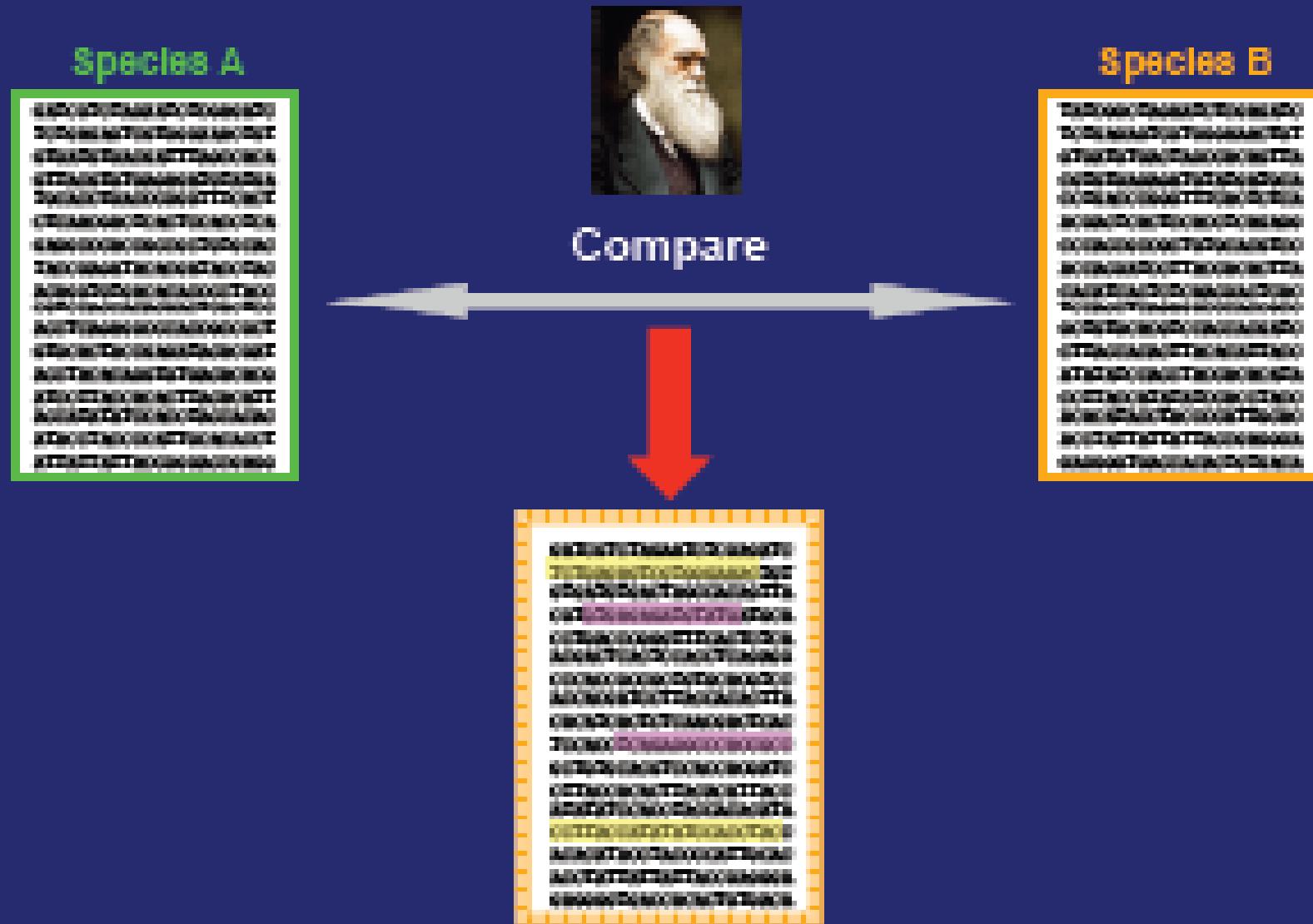
Chromosomal Functional Elements

Undiscovered Functional Elements (NOT Yet in Textbooks!)

Poor Inventory at Present The Genomic Landscape: *circa 2010, Eric Green*

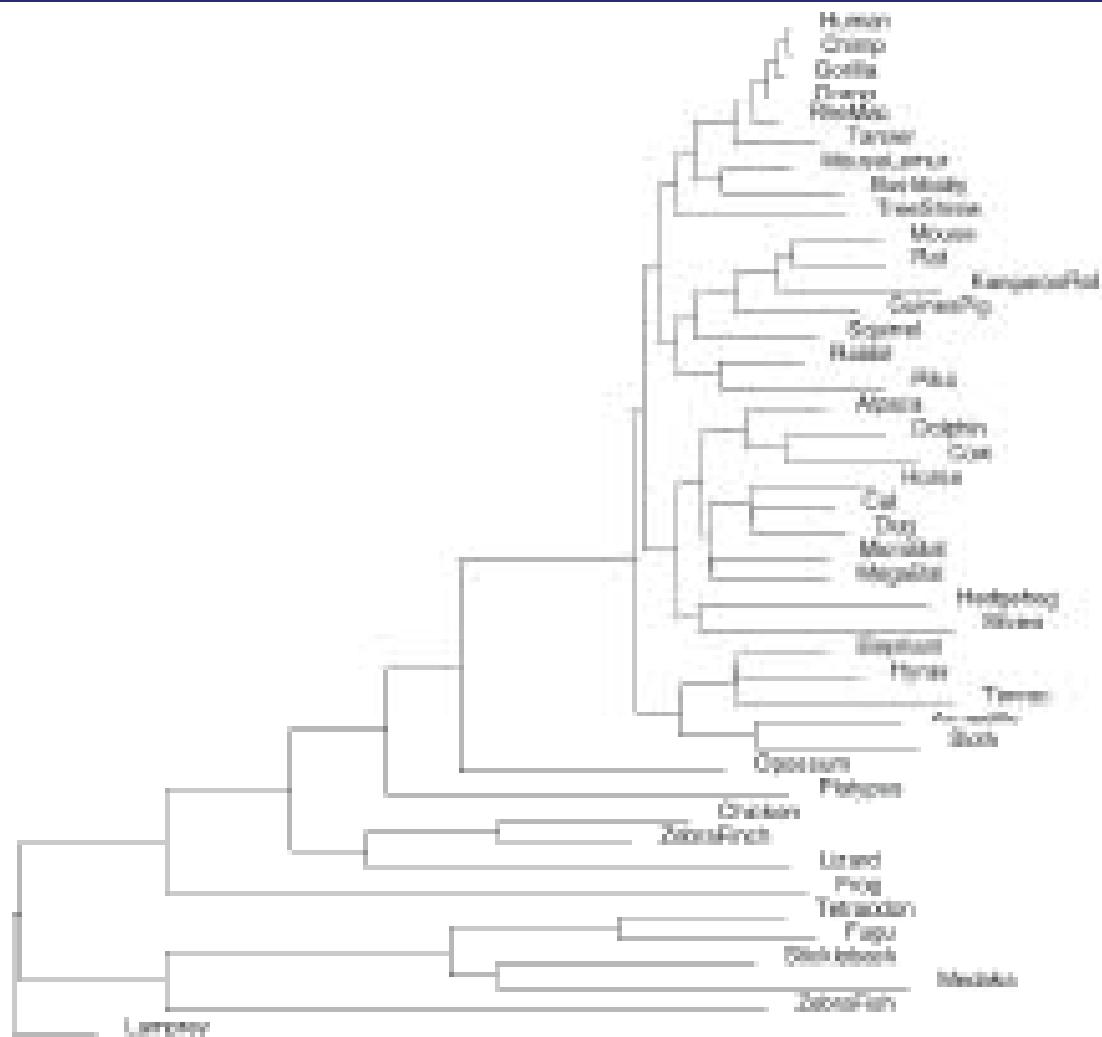
Comparative Sequence Analysis

Using the 'Experiments of Evolution' to Decode the Human Genome



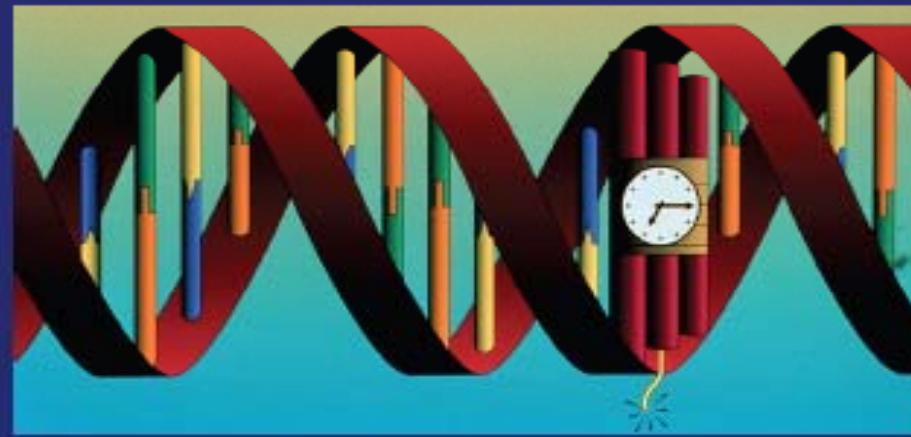
Sequences in Common (i.e., ‘Conserved’ or ‘Constrained’) The Genomic Landscape: *circa* 2010, Eric Green

22 Additional Mammalian Genome Sequences (@ Low Redundancy)



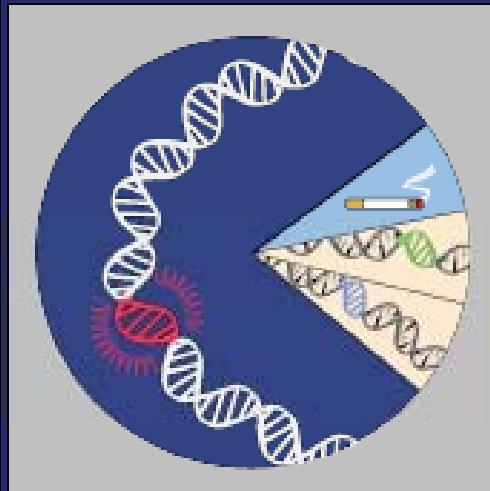
The Genomic Landscape: circa 2010, Eric Green

All humans are ~99.7% identical at the DNA sequence level, and yet...

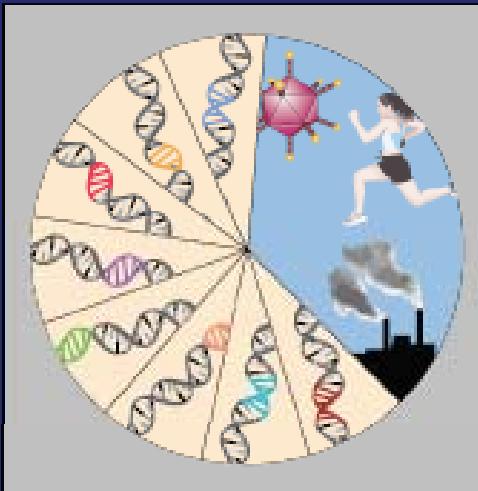


all of us carry a significant number of 'glitches' in our genomes.

Genomic Architecture of Genetic Diseases



Rare, Simple, Monogenic,
Mendelian...



Common, Complex, Multigenic,
Non-Mendelian...

The Genomic Landscape: *circa 2010, Eric Green*



International HapMap Project

[Home](#) | [About the Project](#) | [Data](#) | [Publications](#) | [Tutorial](#)



THE HUMAN VARIOME PROJECT

sharing data • reducing disease

[Home](#) [About](#) [Activities](#) [Publications](#) [Recommendations](#) [Meetings](#) [News](#) [Links](#) [Blogs](#) [Members](#)

Nuevas tecnologías de secuenciación de ADN

Eric D. Green, M.D., Ph.D.

Techniques for Genome Mapping & Sequencing

Genome sequencing in microfabricated high-density picolitre reactors

Margulies M et al. (2005)

Toward the \$1000 human genome

Bennett (2004)

Accurate Multiplex Polony Sequencing of an Evolved Bacterial Genome

Shendure et al. (2005)

Bead-based sequencing by ligation

Bennett et al. (2005)

Perspective

Emerging technologies in DNA sequencing

Michael L. Metzker

Human Genome Sequencing Center and Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, Texas 77030, USA

Demand for DNA sequence information has never been greater, yet current Sanger technology is too costly, time consuming, and labor intensive to meet this ongoing demand. Applications span numerous research interests, including sequence variation studies, comparative genomics and evolution, forensic, and diagnostic and applied therapeutics. Several emerging technologies show promise of delivering next-generation solution for fast and affordable genome sequencing. In this review article, the DNA polymerase-dependent strategies of Sanger sequencing, single nucleotide addition, and cyclic reversible termination are discussed, along with potential challenges these technologies face in their development for ultrafast sequencing.

Metzker (2005)

SCIENTIFIC AMERICAN

Know Your DNA

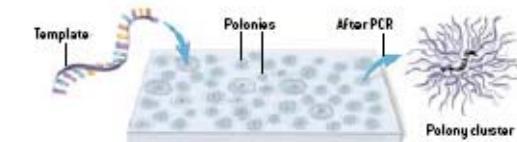
Inexpensive gene readers will soon unlock the secrets in your personal double helix

Alternatives to Toxic Tests on Animals

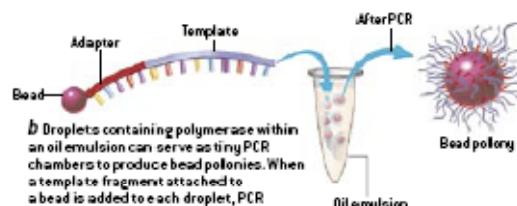
Church (2006)

AMPLIFICATION

Because light signals are difficult to detect at the scale of a single DNA molecule, base-extension or ligation reactions are often performed on millions of copies of the same template strand simultaneously. Cell-free methods (a and b) for making these copies involve PCR on a miniaturized scale.

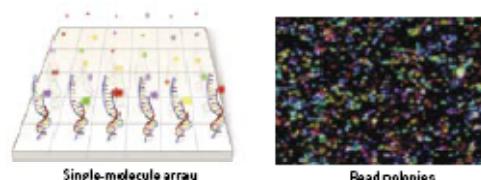


a Polonies—polymerase colonies—created directly on the surface of a slide or gel each contain a primer, which a template fragment can find and bind to. PCR within each polony produces a cluster containing millions of template copies.



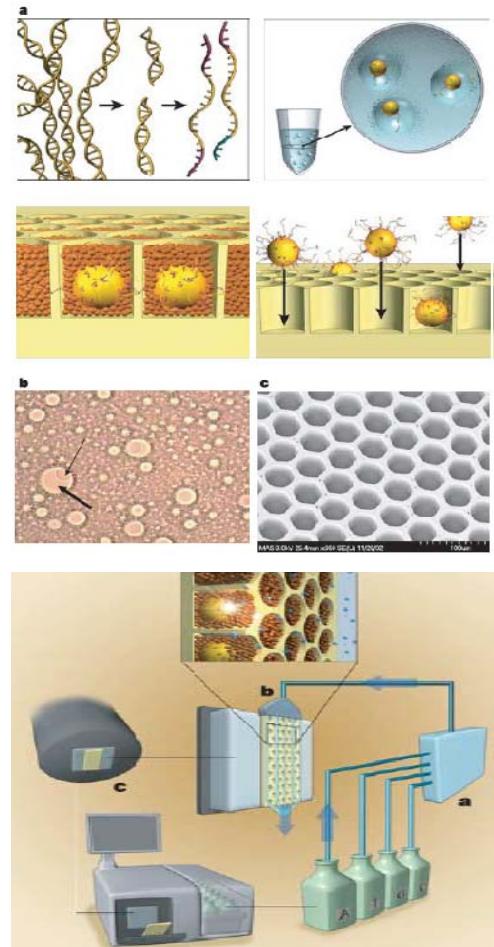
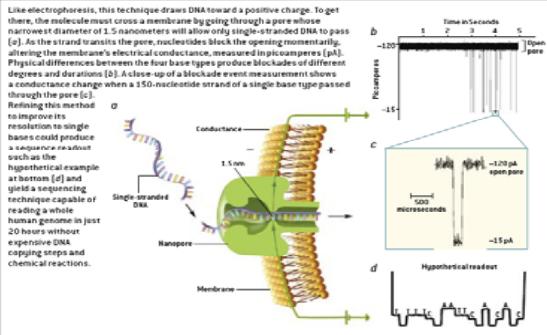
MULTIPLEXING

Sequencing thousands or millions of template fragments in parallel maximizes speed. A single-molecule base-extension system using fluorescent signal detection, for example, places hundreds of millions of different template fragments on a single array (below left). Another method immobilizes millions of bead polonies on a gel surface for simultaneous sequencing by ligation with fluorescence signals, shown in the image at right below, which represents 0.01 percent of the total slide area.



NANOPORE SEQUENCING

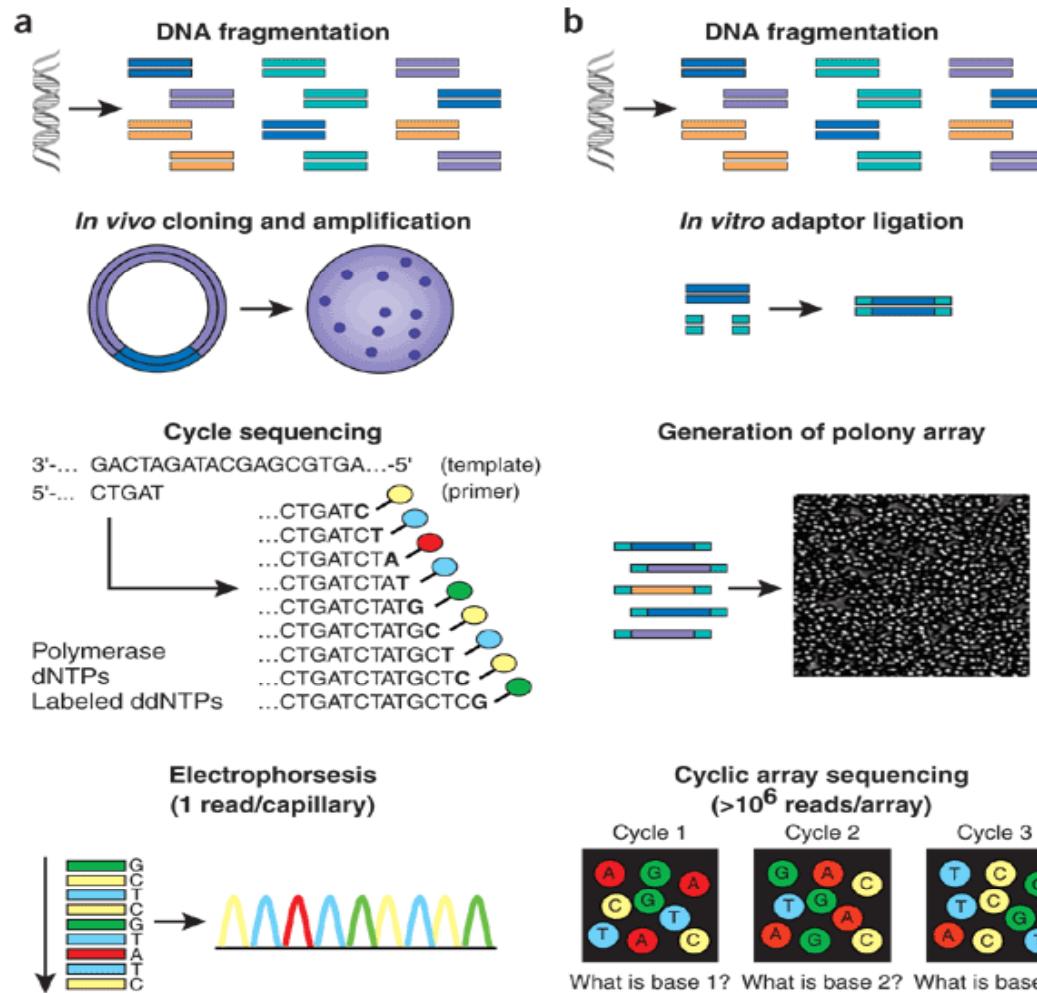
Like electrophoresis, this technique draws DNA toward a positive charge. To get there, the molecule crosses a membrane by going through a pore where the diameter is only about 2.5 nanometers. The pore opens and closes (a). As the strand transits the pore, nucleotides block the opening momentarily, altering the membrane's electrical conductance, measured in picopascals (pA). Physical differences between the four base types produce blocks of different degrees and durations (b). A close-up of a blockade event measurement shows through-the-pore (c) when a 150-nucleotide strand of a single base type passed through the pore (d). Refining this method to improve resolution to single bases could produce a sequence readout such as this hypothetical example at bottom (d) and yield a sequencing technology capable of reading a whole human genome in just 20 hours without expensive DNA copying steps and chemical reactions.



DNA Sequencing Technologies

Method	Feasibility	Read Length	Data Quality	Raw Data Production
Sanger Sequencing	Well Established	Long (800-1200 bases)	+++	+
Stepwise Synthesis	Becoming Established	Short (25-100 bases)	+	++++
Single Molecule	Far from Established	Long (>1000 bases ?)	???	++++?

Nueva generación de tecnologías de secuenciación



(a) With high-throughput shotgun Sanger sequencing, genomic DNA is fragmented, then cloned to a plasmid vector and used to transform *E. coli*. For each sequencing reaction, a single bacterial colony is picked and plasmid DNA isolated. Each cycle sequencing reaction takes place within a microliter-scale volume, generating a ladder of ddNTP-terminated, dye-labeled products, which are subjected to high-resolution electrophoretic separation within one of 96 or 384 capillaries in one run of a sequencing instrument. As fluorescently labeled fragments of discrete sizes pass a detector, the four-channel emission spectrum is used to generate a sequencing trace. (b) In shotgun sequencing with cyclic-array methods, common adaptors are ligated to fragmented genomic DNA, which is then subjected to one of several protocols that results in an array of millions of spatially immobilized PCR colonies or 'polonies'¹⁵. Each polony consists of many copies of a single shotgun library fragment. As all polonies are tethered to a planar array, a single microliter-scale reagent volume (e.g., for primer hybridization and then for enzymatic extension reactions) can be applied to manipulate all array features in parallel. Similarly, imaging-based detection of fluorescent labels incorporated with each extension can be used to acquire sequencing data on all features in parallel. Successive iterations of enzymatic interrogation and imaging are used to build up a contiguous sequencing read for each array.

Table 1 Second-generation DNA sequencing technologies

	Feature generation	Sequencing by synthesis	Cost per megabase	Cost per instrument	Paired ends?	1° error modality	Read-length	References
454	Emulsion PCR	Polymerase (pyrosequencing)	~\$60	\$500,000	Yes	Indel	250 bp	14,20
Solexa	Bridge PCR	Polymerase (reversible terminators)	~\$2	\$430,000	Yes	Subst.	36 bp	17,22
SOLiD	Emulsion PCR	Ligase (octamers with two-base encoding)	~\$2	\$591,000	Yes	Subst.	35 bp	13,26
Polonator	Emulsion PCR	Ligase (nonamers)	~\$1	\$155,000	Yes	Subst.	13 bp	13,20
HeliScope	Single molecule	Polymerase (asynchronous extensions)	~\$1	\$1,350,000	Yes	Del	30 bp	18,30

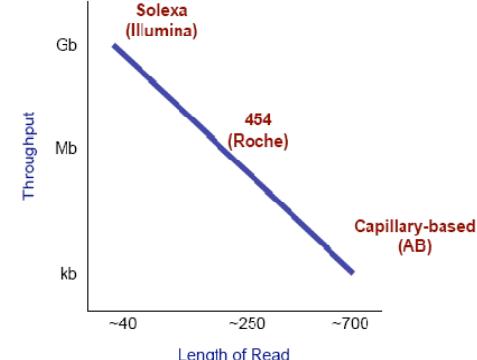
The pace with which the field is moving makes it likely that estimates for costs and read-lengths will be quickly outdated. Vendors including Roche Applied Science, Illumina, and Applied Biosystems have major upgrade releases currently in progress. Estimated costs-per-megabase are approximate and inclusive only of reagents. Read-lengths are for single tags. Subst., substitutions; indel, insertions or deletions; del, deletions.

Nature Biotechnology 26, 1135 - 1145 (2008)

Advances in DNA sequencing technologies

Technology	Approach	Read length	Bp per run	Company name and Web site
Automated Sanger sequencer	Synthesis in the presence of dye terminators	Up to 900 bp	96 kb	Applied Biosystems www.appliedbiosystems.com
ABI3730xl				
454/Roche FLX system	Pyrosequencing on solid support	200–300 bp	80–120 Mb	Roche Applied Science www.roche-applied-science.com
Illumina/Solexa	Sequencing by synthesis with reversible terminators	30–40 bp	1 Gb	Illumina, Inc. http://www.illumina.com/
ABI/SOLiD	Massively parallel sequencing by ligation	35 bp	1–3 Gb	Applied Biosystems www.appliedbiosystems.com

Trade-offs with Newer Sequencing Technologies



R McCombie, 2008. Next Generation Sequencing Symposium. Santa Fe

Genomics 92 (2008) 255–264

Procesamiento de secuencias en proyectos genómicos

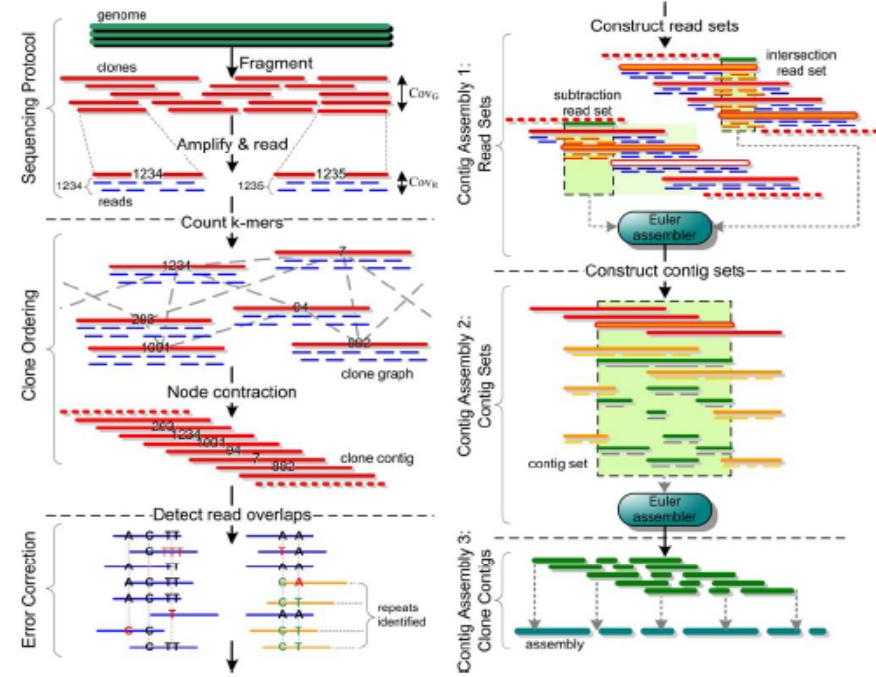
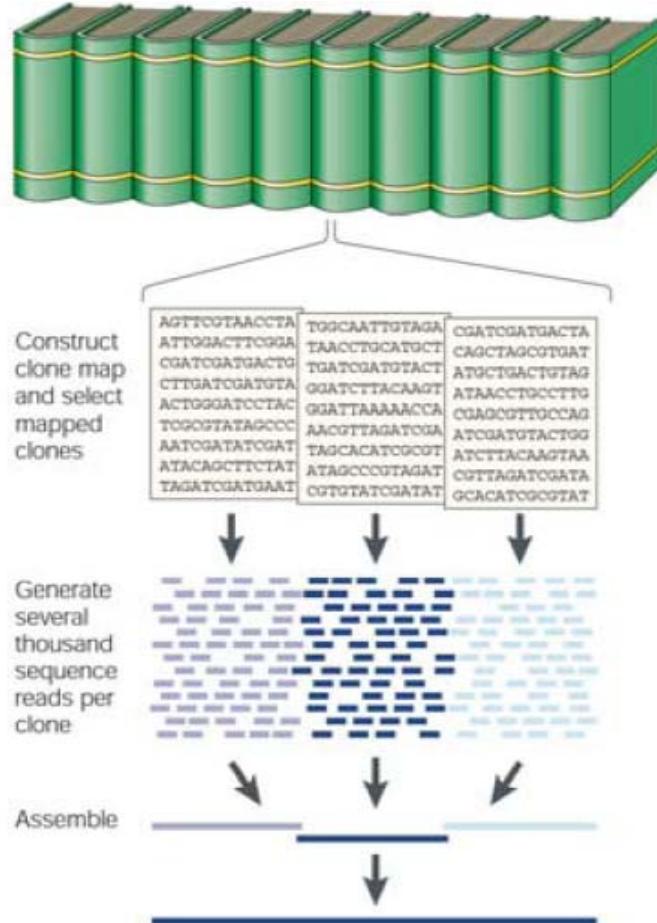
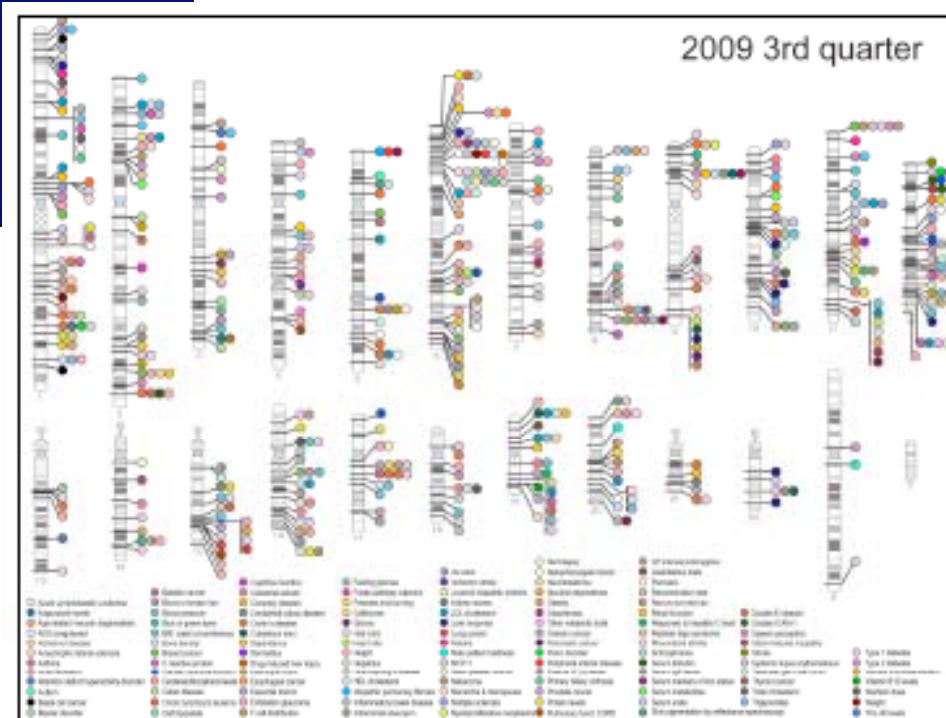
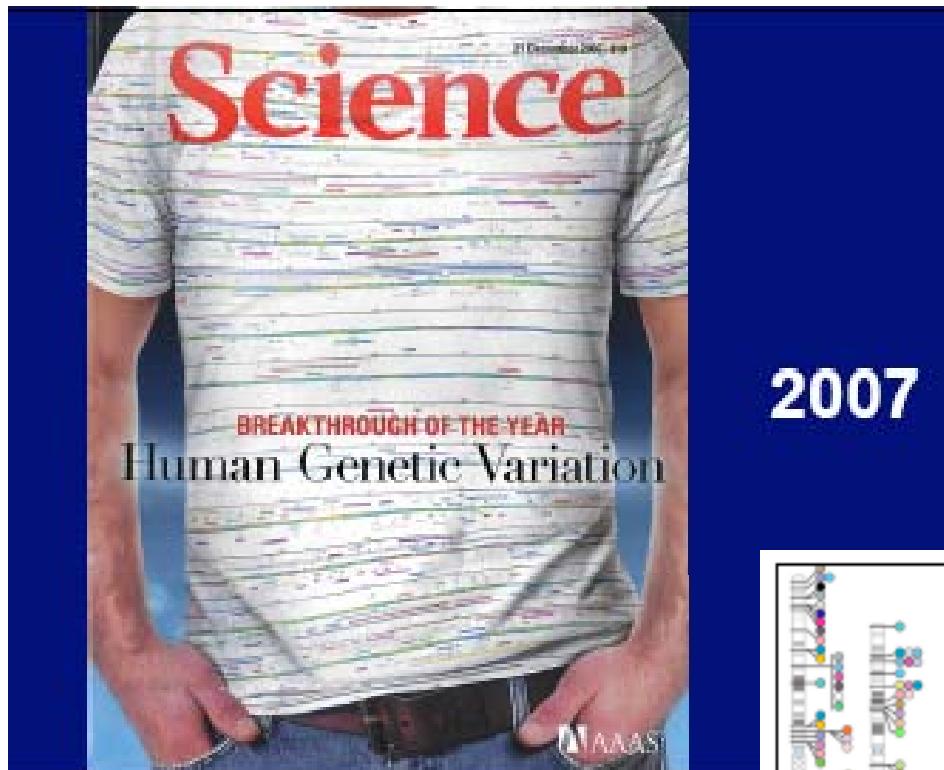


Figure 1. Sequencing protocol and assembly methodology. Reads are obtained in a hierarchical sequencing protocol with high genome-dope coverage and low clone-read coverage. From the k-mer content of each clone we construct a clone graph whose edge weights reflect the likely clone proximities, and from this our clone ordering algorithm determines the clone config. Next, we find all putative read overlaps by only looking in nearby clones and perform error correction. In three stages of contig assembly we 1) create read sets via set operations that consist of reads from multiple overlapping clones within small clone subsections and assemble using Euler; 2) combine contigs resulting from the previous stage in clone-sized contig sets for assembly; and 3) use a scalable assembler to merge entire clone contigs.
doi:10.1371/journal.pone.0000484.g001

PLoS ONE 2(5): e484.
doi:10.1371/journal.pone.0000484



<http://demo.decodeme.com/research-catalog>

Human Genome Sequence

>\$1,000,000,000

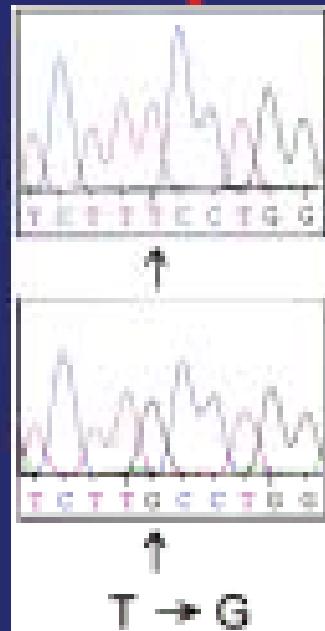
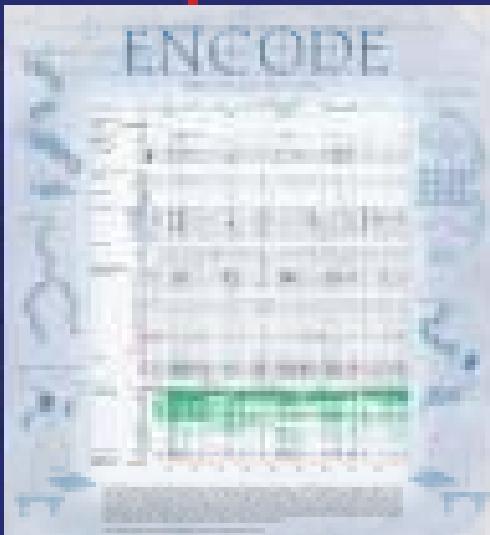
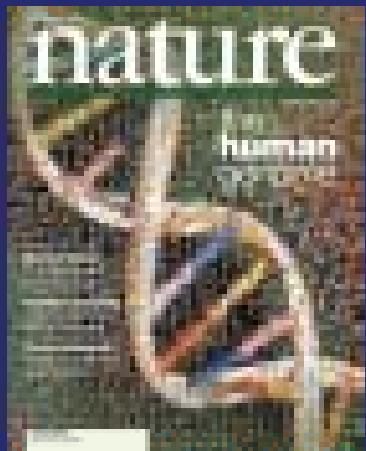


~\$1,000

The Pathway to Genomic Medicine

Interpreting
the Human
Genome Sequence

Implicating
Genetic Variants
with Human Disease



Realization of
Genomic Medicine

The Genomic Landscape: *circa 2010, Eric Green*