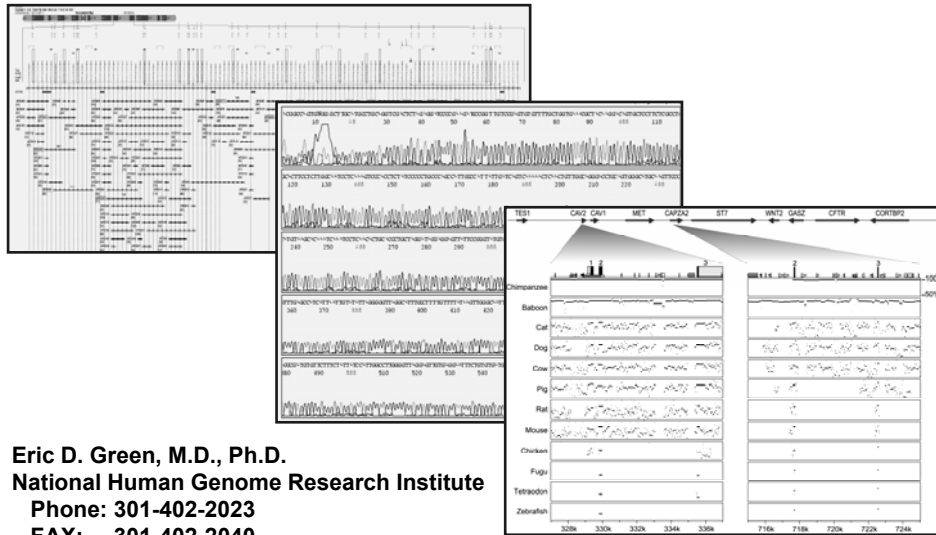


Techniques for Genome Mapping & Sequencing



Eric D. Green, M.D., Ph.D.
 National Human Genome Research Institute
 Phone: 301-402-2023
 FAX: 301-402-2040
 E-Mail: egreen@nhgri.nih.gov

Foundational Milestones in Genetics & Genomics



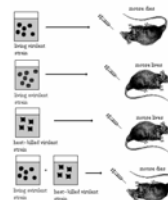
Mendel

1865



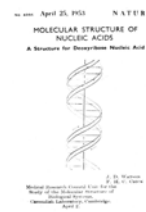
Miescher

1871



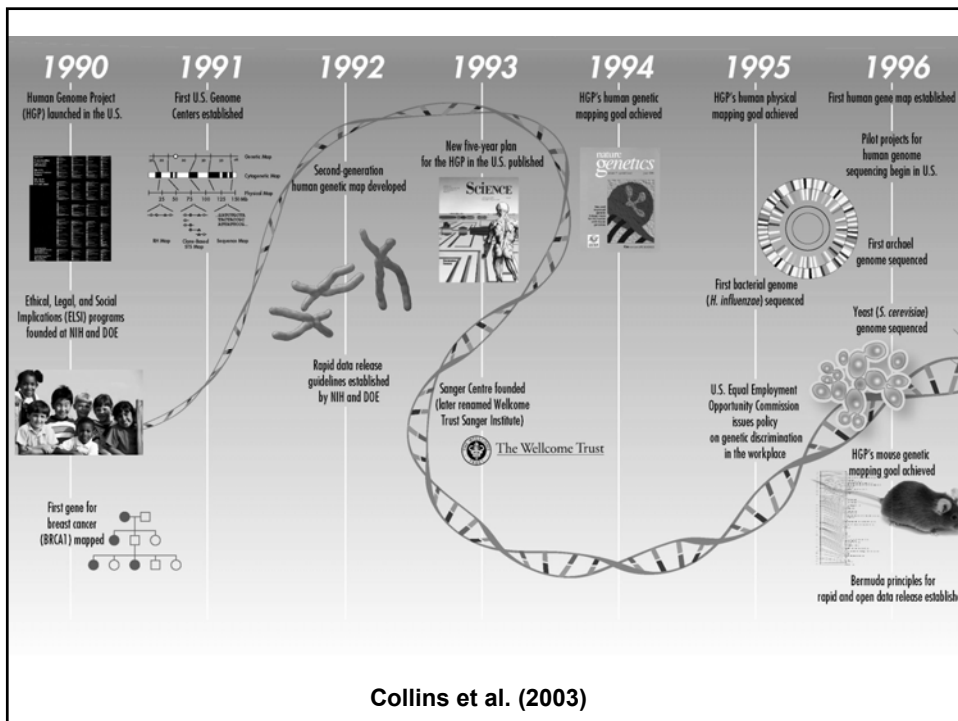
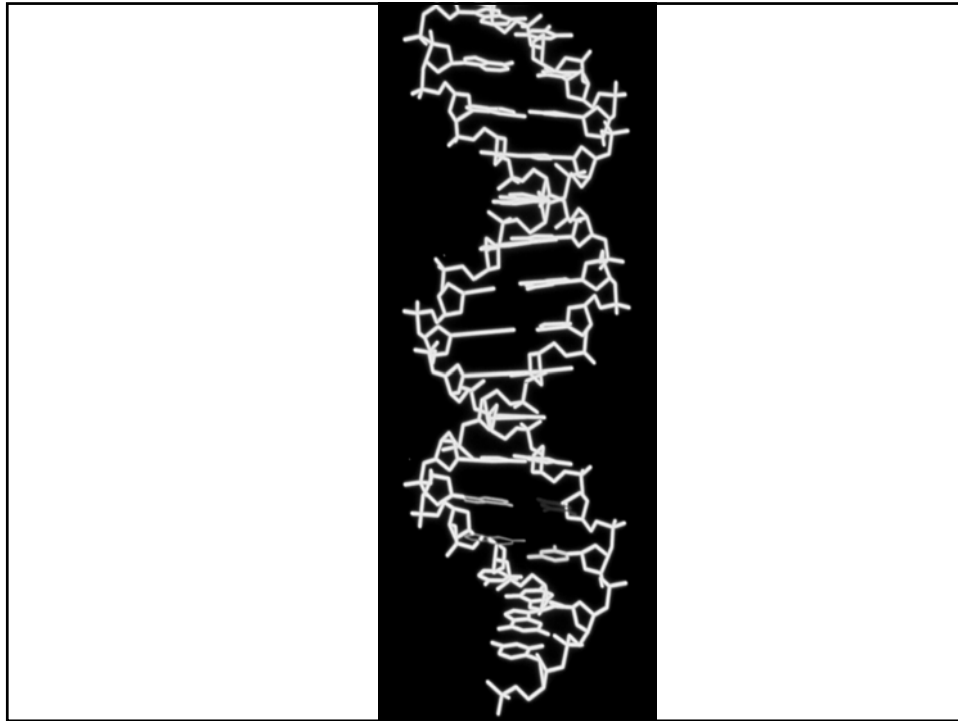
Avery

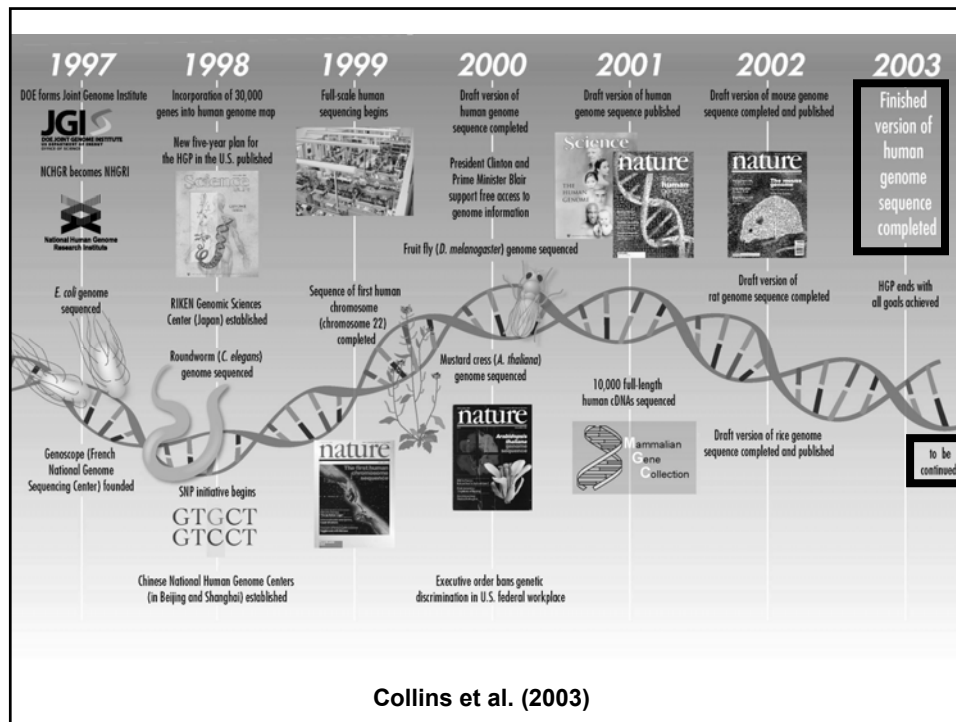
1944



Watson & Crick

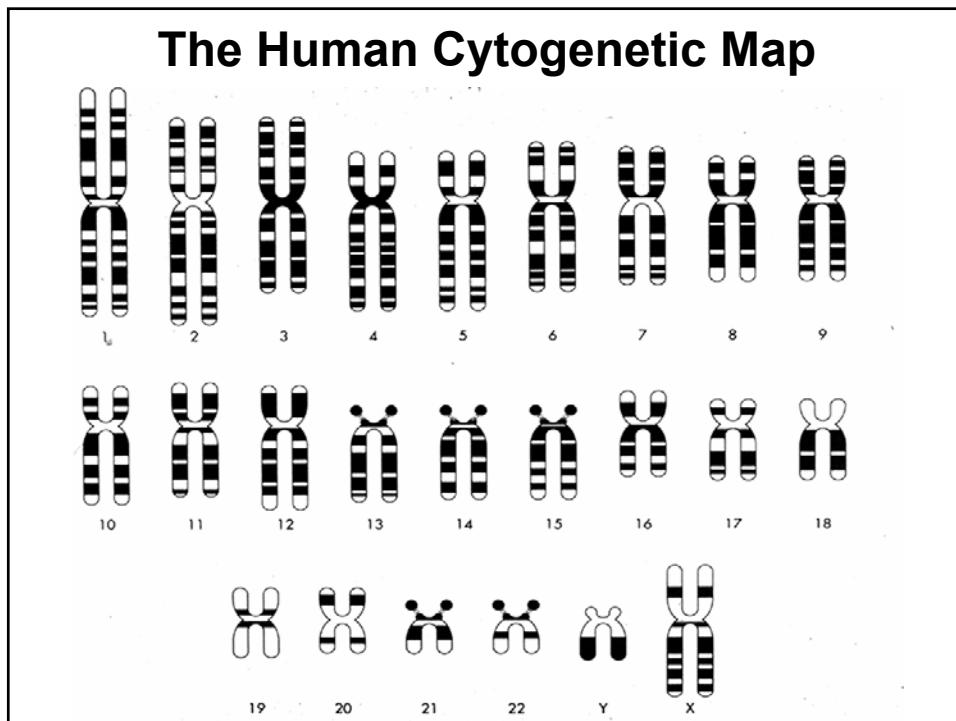
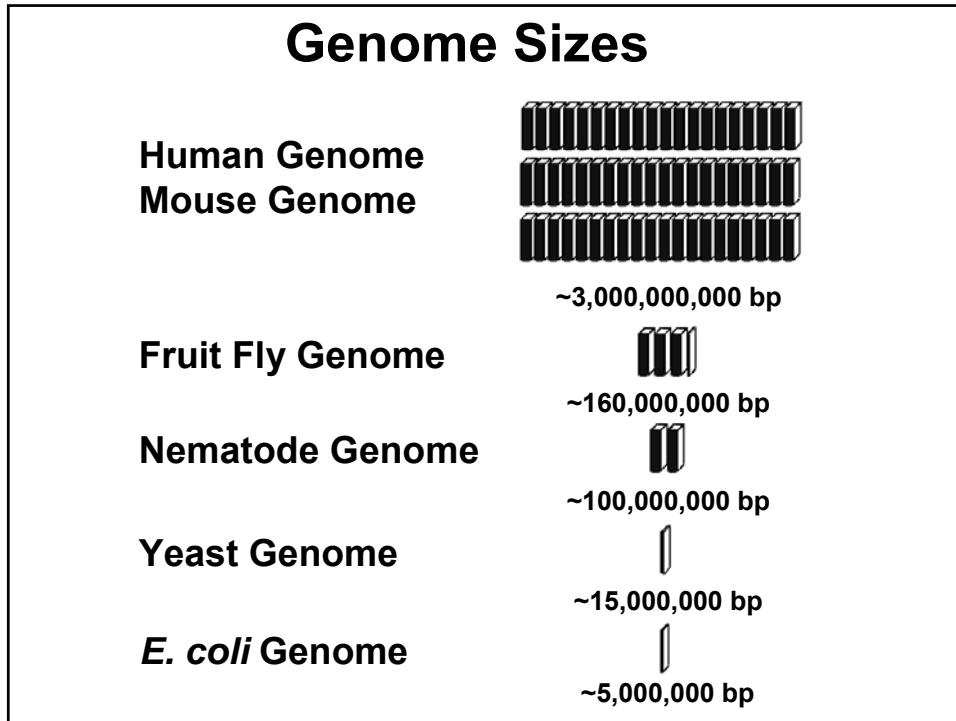
1953

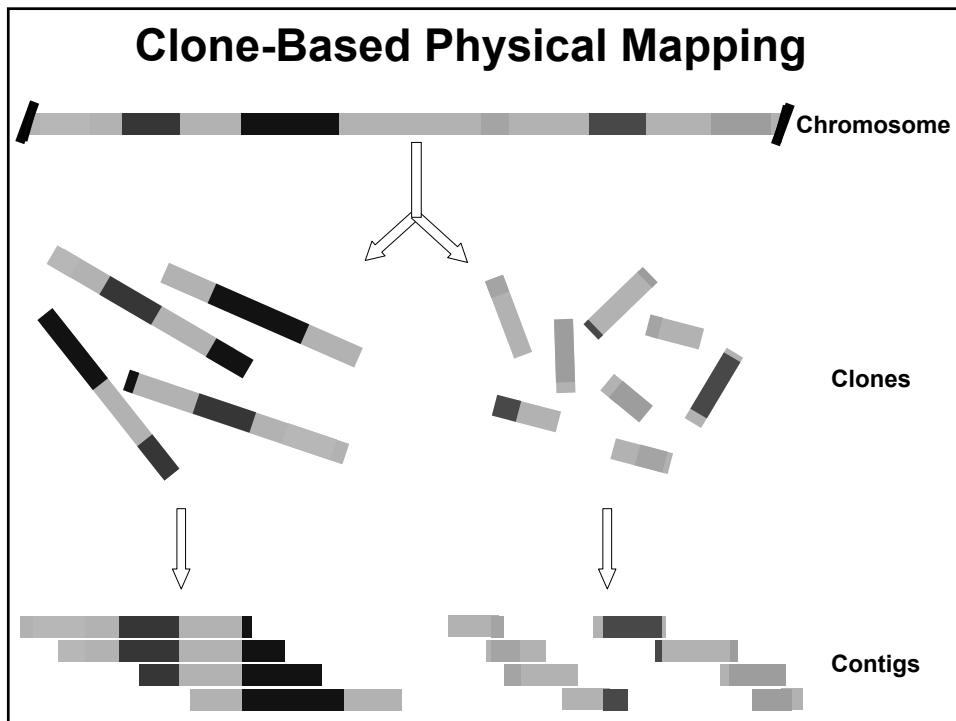
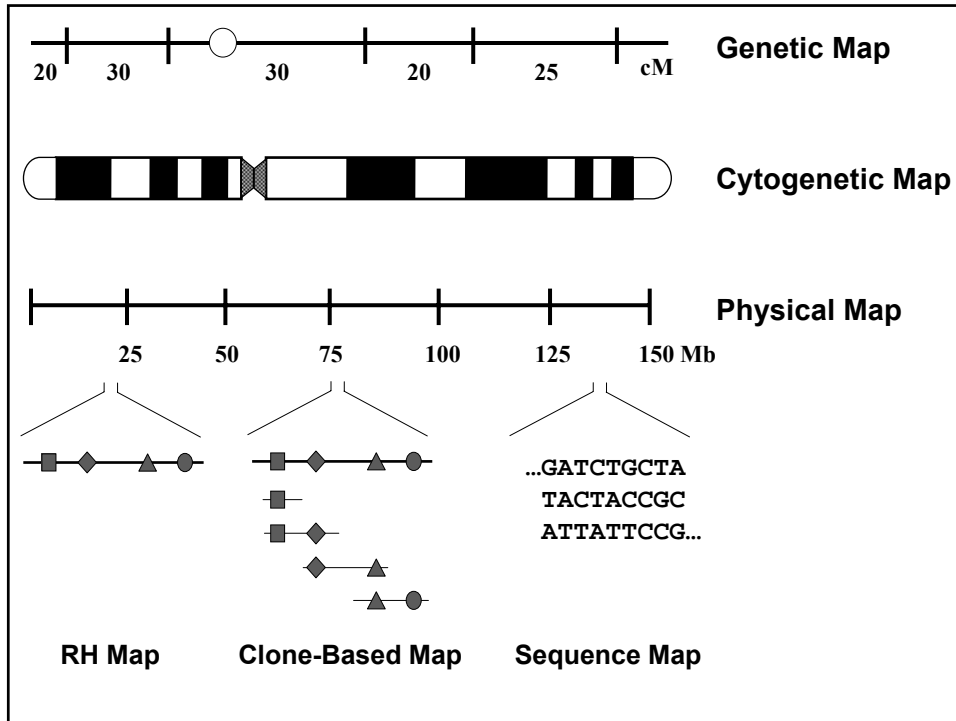


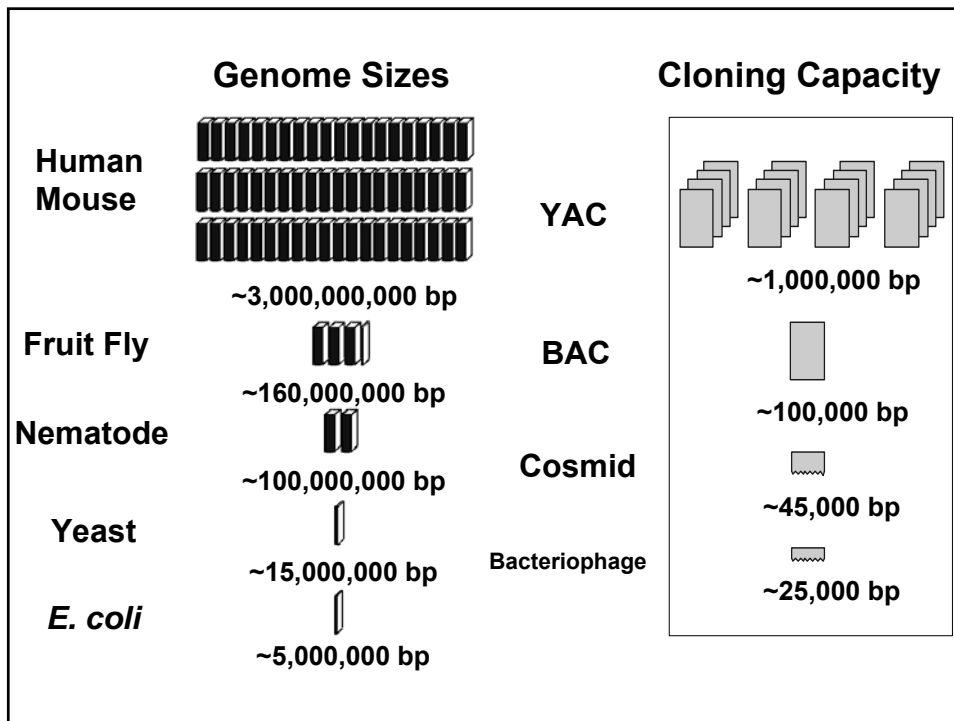
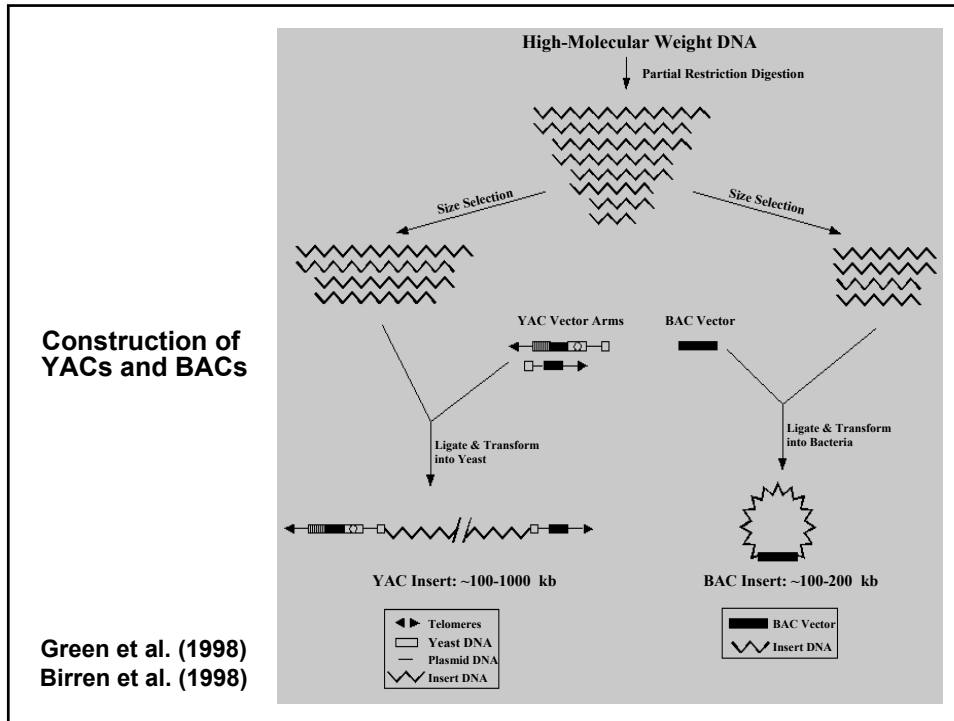


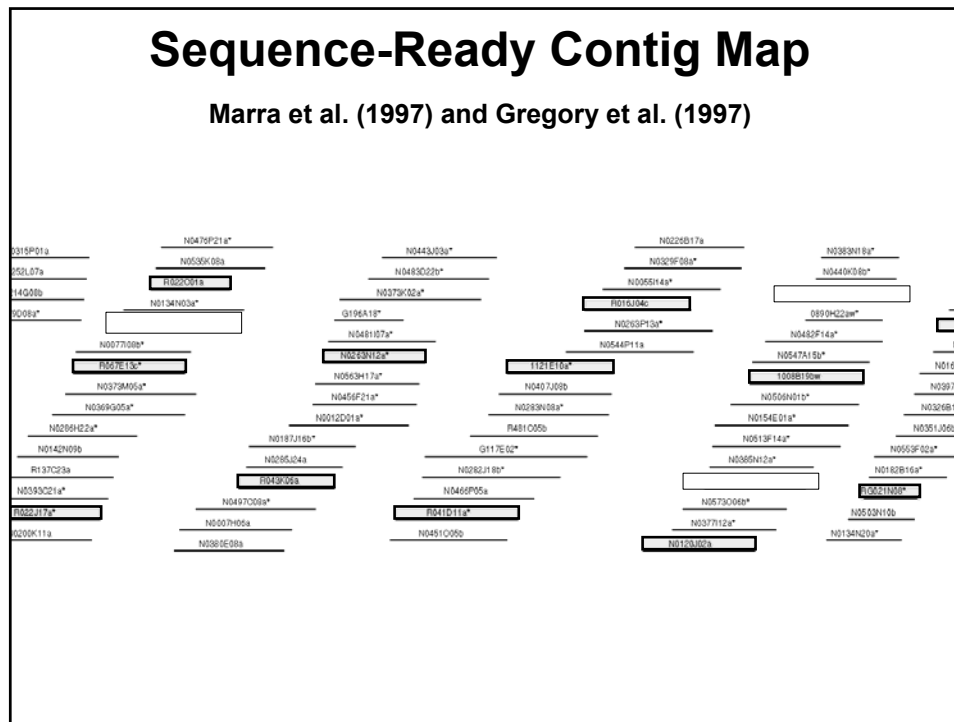
Outline

- I. Fundamentals of Genome Mapping
- II. Fundamentals of Genome Sequencing
- III. Mapping & Sequencing in the Human Genome Project
- IV. Comparative Sequencing
- V. New DNA Sequencing Technologies







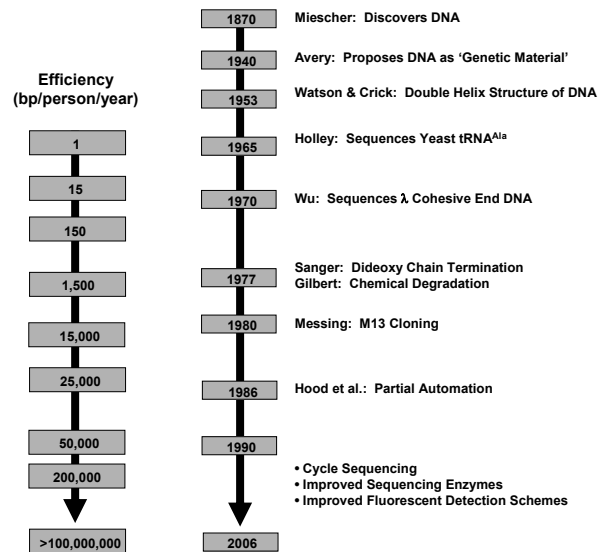


Physical Mapping: Future Prospects

- **Strategies for Physical Mapping have Advanced Greatly in the Sequence-Based Era**
- **Close Interplay of Mapping and Sequencing in the Exploration of Genomes**
- **Availability of Many BAC Libraries is Allowing Physical Mapping of More Species' Genomes**

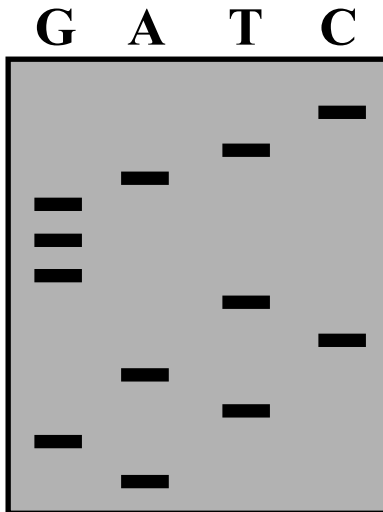
DNA Sequencing

History of DNA Sequencing



Adapted from Messing & Liaca, *PNAS* (1998)

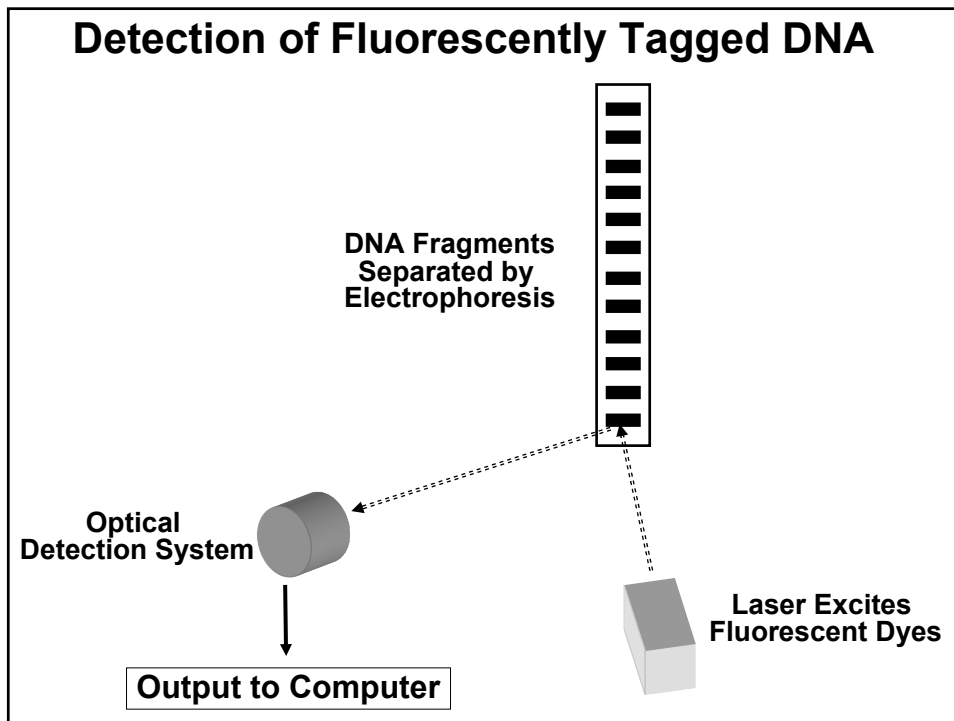
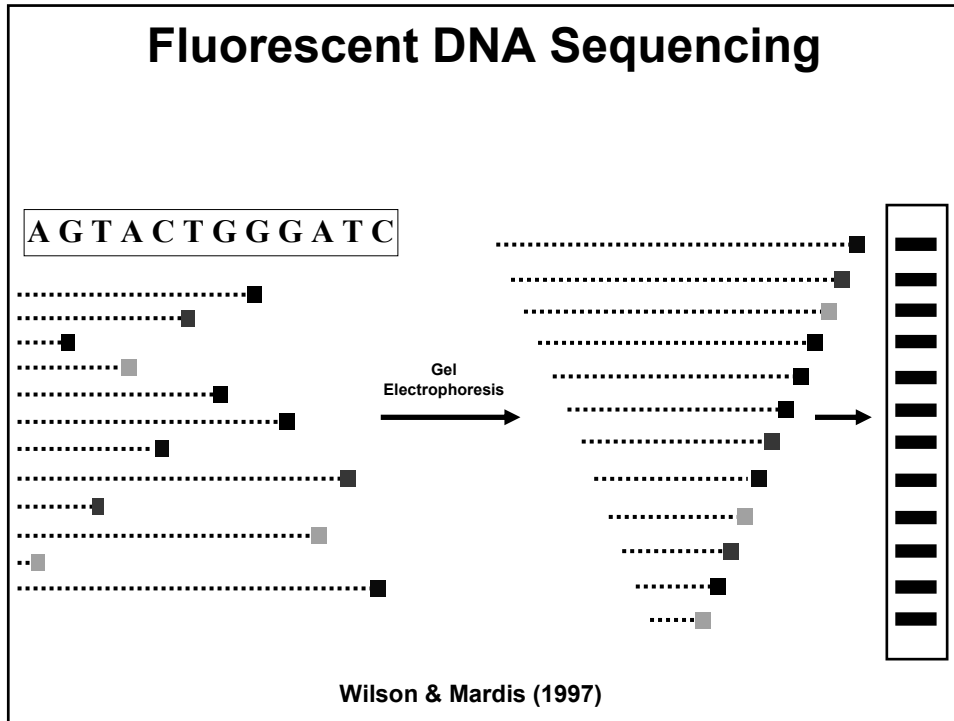
DNA Tagged with Radioactivity

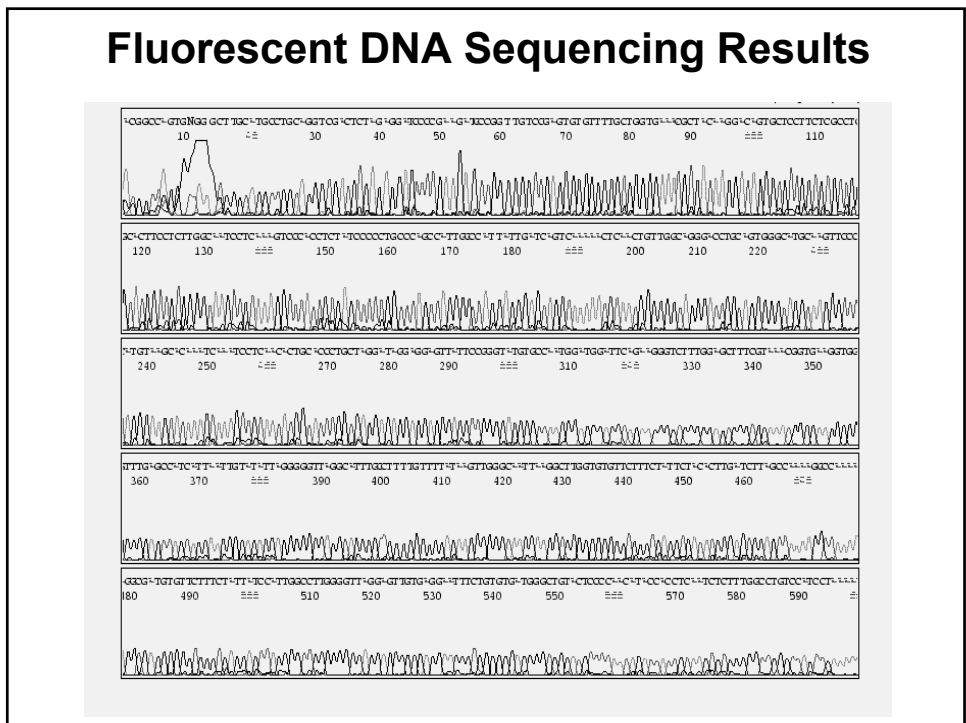
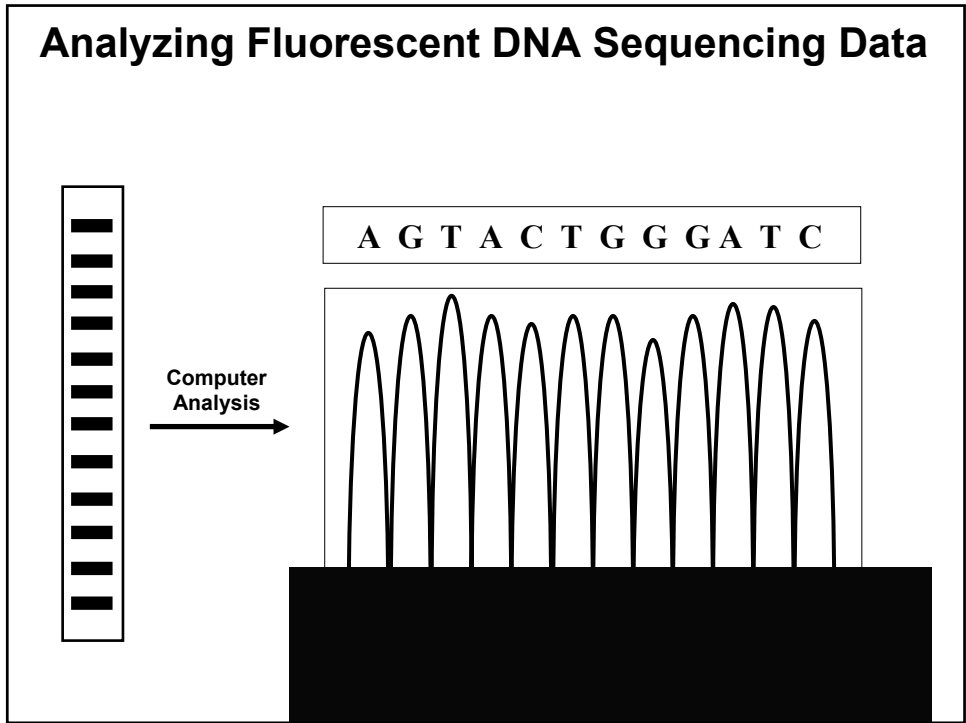


G: G Reaction
A: A Reaction
T: T Reaction
C: C Reaction

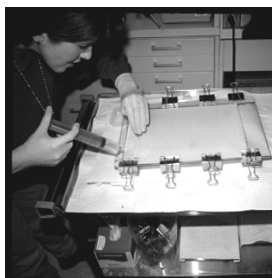
Radioactive Sequencing







Slab Gel-Based DNA Sequencing Instruments

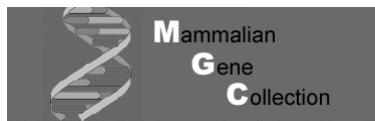


Capillary-Based DNA Sequencing Instruments



Large-Scale cDNA Sequencing

- **ESTs:** Expressed-Sequence Tags
- **SAGE:** Serial Analysis of Gene Expression
- **Full-Insert (Full-Length) cDNA Sequencing**



mgc.nci.nih.gov

Gerhard et al. (2004)

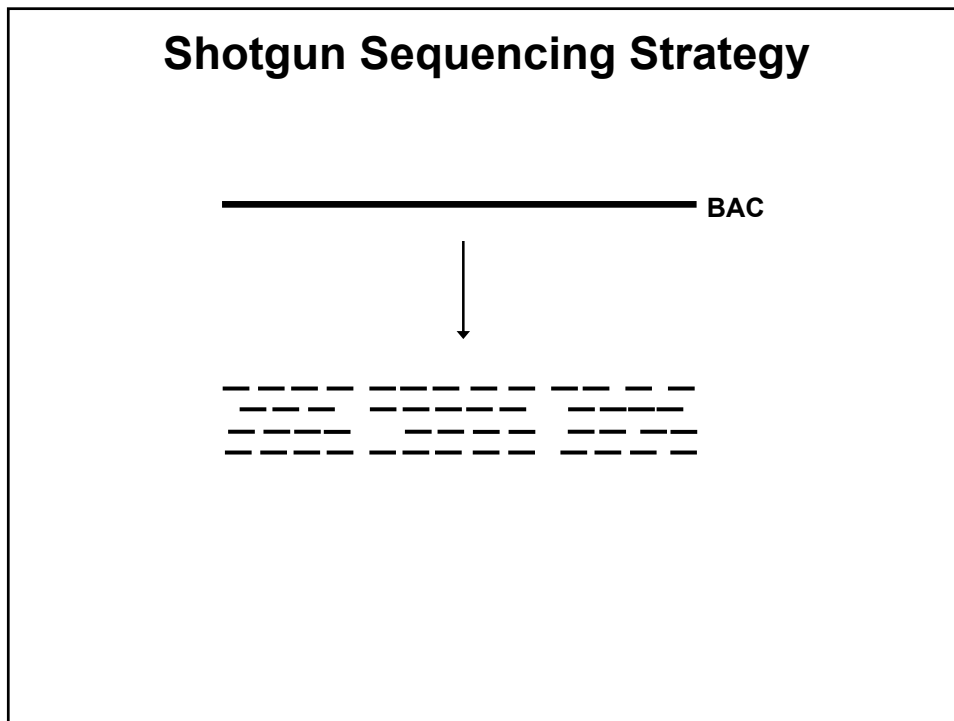
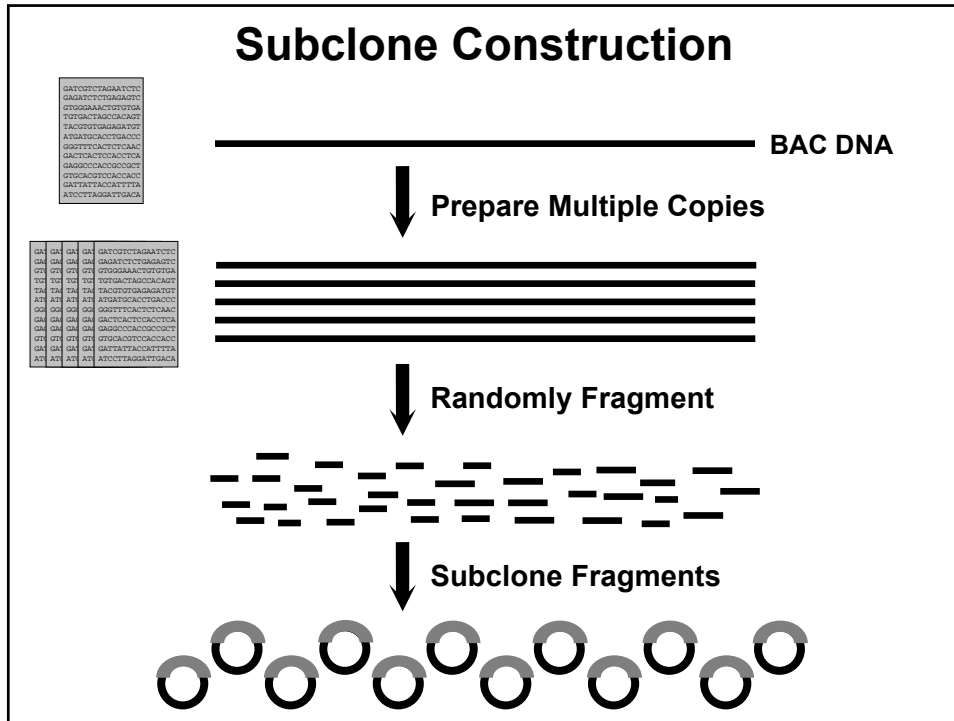
Large-Scale Genome Sequencing

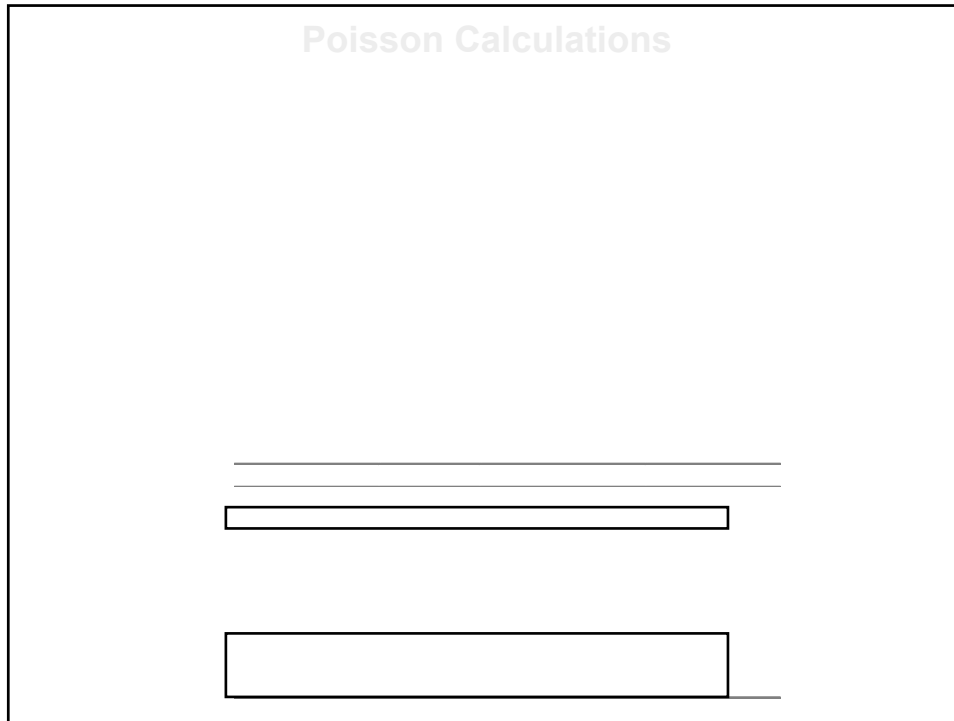


Shotgun Sequencing

Wilson & Mardis (1997)

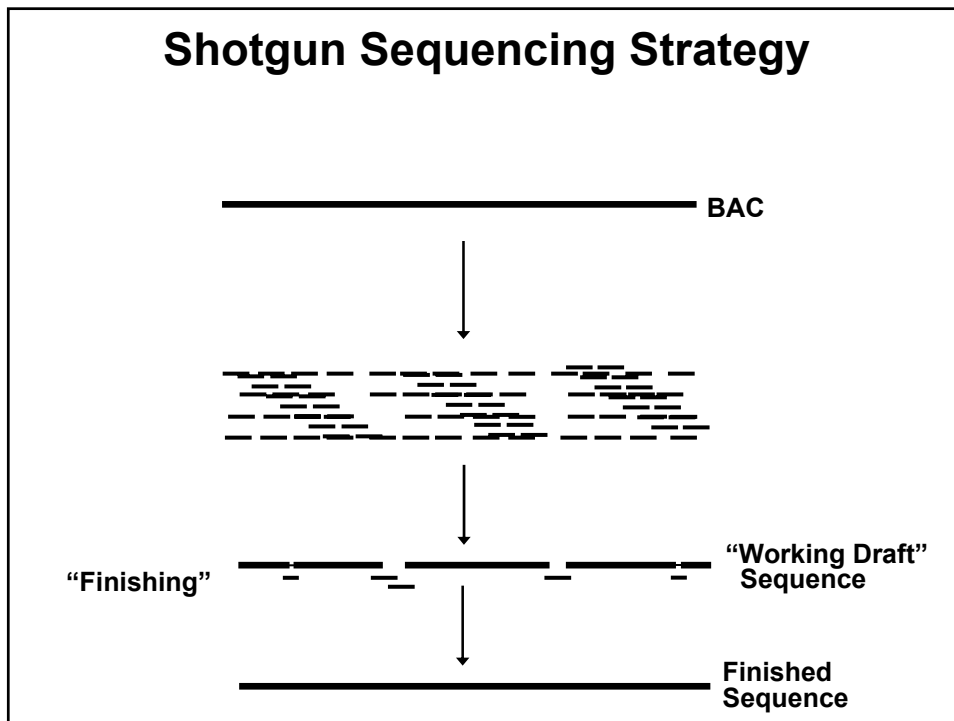
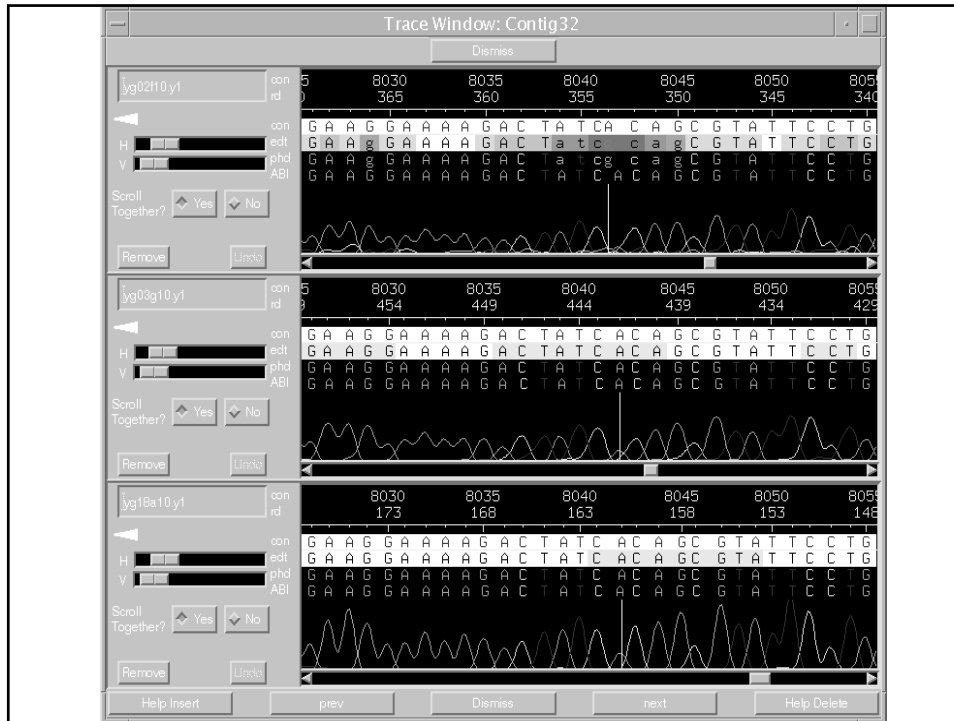
Green (2001)



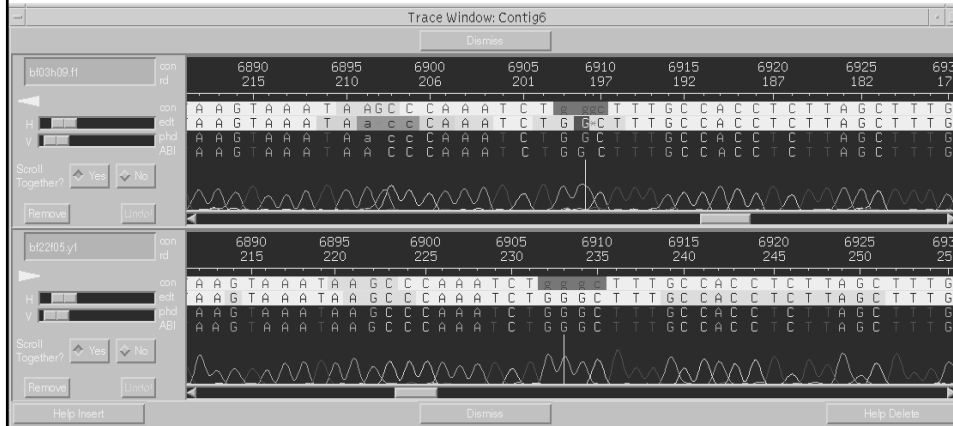


Shotgun Sequence Assembly

"Consed" (Gordon et al., 1998)



Sequence Finishing: Resolving Ambiguities



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

Historically Significant Genome Sequencing Projects

Microbial Genome Sequences

TIGR-CMR
 Click here to take the CMR user feedback survey
 Home|Genomes|Searches|Comparative Analyses|Gene Lists|Carts|Downloads

Welcome to the Comprehensive Microbial Resource (CMR) Home Page

The Comprehensive Microbial Resource (CMR) is a free website used to display information on all of the publicly available, complete prokaryotic genomes. In addition to the convenience of having all of the organisms on a single website, common data types across all genomes in the CMR make searches more meaningful, and cross genome analysis highlight differences and similarities between the genomes. [More Information] [Publication Information]

NEW on the CMR

August 11, 2006: CMR Feedback Survey: TIGR is currently seeking funding for the maintenance and expansion of the CMR. Please help us with our grant application by taking our user feedback survey.

For information on the latest CMR updates, subscribe to the CMR Mailing List.

Prokaryotic Annotation and Analysis Classes: June 13-15, August 8-10, October 10-12, 2006

Visit a CMR page for an individual genome

Acidithiobacillus ferrooxidans ATCC 23270

Acinetobacter sp. ADP1

Actinomyces naeslundii MG1

Aeropyrum pernix K1

Agrobacterium tumefaciens C58 Cereon

Agrobacterium tumefaciens C58 UWash

Anabaena variabilis ATCC 29413

Anaplasma marginale St Maries

CMR Genomes: Data Release 19.0

	Complete	Incomplete	Totals
Bacteria	279	17	296
Archaea	23	0	23
Viruses	3	0	3
Totals	305	17	322

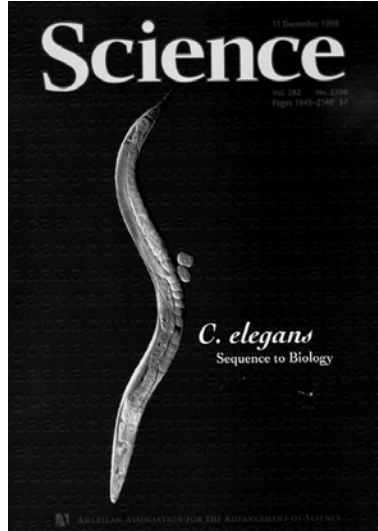
www.tigr.org

First Eukaryotic Genome Sequence

The yeast genome directory

Goffeau et al. (1997)

First Animal Genome Sequence

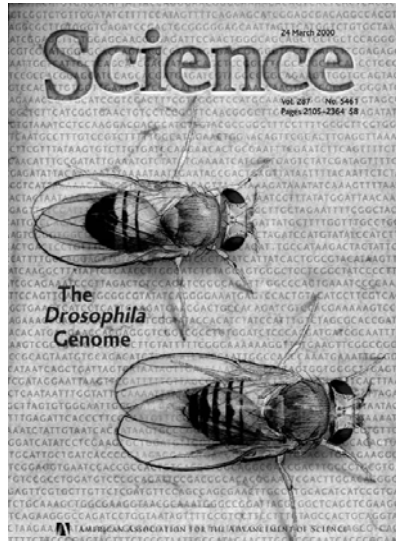


Genome Sequence of the Nematode *C. elegans*: A Platform for Investigating Biology

The *C. elegans* Sequencing Consortium*

C. elegans Sequencing Consortium (1998)

Second Animal Genome Sequence



The Genome Sequence of *Drosophila melanogaster*

Mark D. Adams,^{1*} Susan E. Celniker,² Robert A. Holt,¹ Cheryl A. Evans,¹ Jeannine D. Gocayne,¹ Peter G. Amanatides,¹ Steven E. Scherer,¹ Peter W. Li,¹ Roger A. Hoskins,¹ Richard F. Gallie,¹ Reed A. George,² Suzana E. Lewis,⁴ Stephen Richards,¹ Michael Ashburner,¹ Scott W. Henderson,¹ Granger G. Sutton,¹ Jennifer R. Wortman,¹ Mark D. Vandell,¹ Qing Zhang,¹ Lin X. Chen,¹ Rhonda C. Brandon,¹ Yu-Hui C. Rogers,¹ Robert G. Blazaj,¹ Mark Champe,² Barrett D. Pfeiffer,¹ Kenneth H. Wan,² Clare Doyle,¹ Evan G. Baxter,¹ Gregg Heist,⁴ Catherine R. Nelson,¹ George L. Gabor Miklos,¹ Joseph F. Alari,⁴ Anna Aghayani,¹ Hai-Jin An,¹ Cynthia Andrews-Pfannkoch,¹ Daniela Baldwin,¹ Richard M. Ballieu,¹ Anand Basu,¹ James Baxterdale,¹ Leyla Bayraktaroglu,¹ Ellen M. Beasley,¹ Karen Y. Beeson,¹ P. V. Benos,¹⁸ Benjamin P. Berman,² Deepali Bhandari,¹ Slava Bolshakov,¹¹ Dana Borokova,¹¹ Michael R. Botchan,¹ John Bouch,² Peter Brokstein,¹ Philippe Brotier,¹⁴ Kenneth C. Burtis,¹⁵ Dana A. Busam,¹ Heather Butler,¹⁶ Edward Cadieu,¹⁷ Angela Center,¹ Ishwar Chandra,¹ J. Michael Cherry,¹⁸ Simon Cawley,¹⁸ Carl Dahlke,¹ Lionel B. Davenport,¹ Peter Davies,¹ Beatriz de Pablos,²⁰ Arthur Delcher,² Zuoming Deng,¹ Anne Deslattes Mays,¹ Ian Dew,¹ Suzanne M. Dietsi,¹ Kristina Dodson,¹ Lisa E. Doup,¹ Michael Dumas,²¹ Shannon Dugan-Rocha,¹ Boris C. Dunkov,²² Patrick Dunn,¹ Kenneth J. Durbin,³ Carlos C. Evangelista,¹ Concepcion Ferraz,²³ Steven Ferreira,¹ Wolfgang Fleischmann,¹ Carl Foeiser,¹ Andrei E. Gabrielian,¹ Neha S. Garg,¹ William M. Galbraith,² Ken Glasser,¹ Anna Glödeke,¹ Fangcheng Gong,¹ J. Harley Gorrell,¹ Zhiping Guo,¹ Ping Guan,¹ Michael Harris,¹ Nomi L. Harris,¹ Damon Harvey,¹ Thomas J. Heiman,¹ Judith H. Hernandez,¹ Jarrett Houck,¹ Damon Houston,¹ Kathryn A. Houston,¹ Timothy J. Howland,¹ Hing-Hui Wei,¹ Chinyere Ibegwam,¹ Mena Jalali,¹ Francis Kalush,¹ Gary H. Karpen,²⁴ Zhaod Ke,¹ James A. Kennison,²⁵ Karen A. Ketchum,¹ Bruce E. Kimmel,² Chinnappa D. Kodira,¹ Cheryl Kraft,¹ Saul Kravitz,¹ David Kulp,¹ Zhongwei Lai,¹ Paul Lasko,²⁶ Yixiang Lai,¹ Alexander A. Lavinsky,¹ Jinyin Li,¹ Zhanya Li,¹ Yong Liang,¹ Xiaoying Lin,²⁶ Xiangjun Liu,¹ Bettina Mattali,¹ Tina C. McIntosh,¹ Michael P. McLeod,² Duncan McPherson,¹ Genady Merkulov,¹ Natalia V. Milshina,¹ Clark Mobarry,¹ Joe Morris,¹ Ali Moshrefi,¹ Stephen M. Mount,²⁷ Mea Moy,¹ Brian Murphy,¹ Lee Murphy,²⁸ Donna M. Murray,² David L. Nelson,¹ David R. Nelson,²⁹ Keith A. Nelson,¹ Katherine Niccox,² Deborah R. Nuskern,¹ Joanne M. Paclab,² Michael Palazzolo,² Gjang S. Pittman,³ Sue Pan,¹ John Pollard,¹ Vinita Puri,¹ Martin G. Reese,⁴ Knut Reinert,¹ Karin Remington,¹ Robert D. C. Saunders,²⁸ Frederick Scheeler,¹ Hua Shen,¹ Bixiang Christopher Shue,¹ Inga Sidén-Kiamos,¹ Michael Simpson,¹ Marian P. Skupski,¹ Tom Smith,¹ Eugene Spier,¹ Allan C. Spradling,³⁰ Mark Stapleton,² Renee Strong,¹ Eric Sun,¹ Robert Svoboda,²⁸ Cyndee Tector,¹ Russell Turner,¹ Eli Venter,¹ Alhui H. Wang,¹ Xin Wang,¹ Zhen-Yuan Wang,¹ David A. Wassarman,³¹ George M. Weinstock,¹ Jean Weissenbach,¹ Sherita M. Williams,¹ Trevor Woodage,¹ Kim C. Wortley,¹ David Wu,¹ Song Yang,¹ Q. Allison Yao,¹ Jumei Ye,¹ Ho-Fang Yeh,¹⁸ Jayshree S. Zaveri,¹ Ming Zhao,¹ Guoqiang Zhang,¹ Qi Zhao,¹ Linsheng Zheng,¹ Xiangqun H. Zheng,¹ Fei N. Zhong,¹ Wenyan Zhong,¹ Xiaojun Zhou,¹ Shaoqing Zhu,¹ Xiaohong Zhu,¹ Hamilton O. Smith,¹ Richard A. Gibbs,¹ Eugene W. Myers,¹ Gerald H. Rubin,³² J. Craig Venter¹

Adams et al. (2000)

Human Genome Sequencing Centers



Clone-Based Shotgun Sequencing



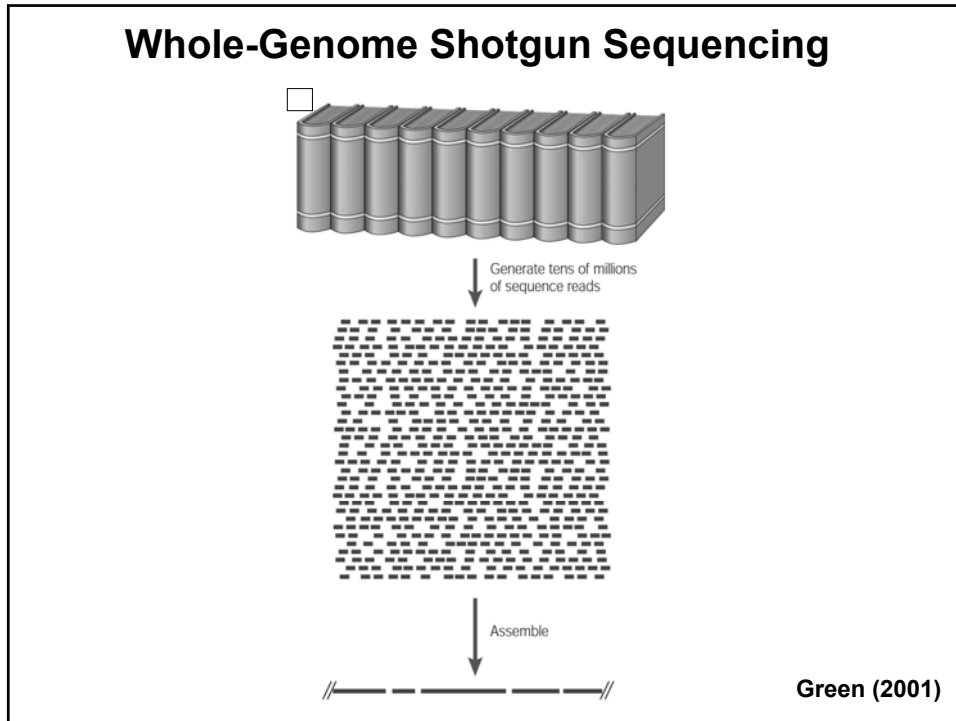
Construct clone map and select mapped clones

AATTCGTAACCTA	TGGCAATTGTAGA	CGATCGATGACTA
ATTGGACTTGGG	TAACTTCAAGCT	CAGTAGCGTAT
CGATCGATGACTG	TGATCGATGACT	ATGCTACTGTAG
CTTGATCGATGTA	GGATCTTCAAGT	ATAGCTCTCTTG
ACTGGATCTTAC	GGATTAAGAACCA	CGAGCTTCCAG
TGCGTATAGCCC	AACOTTAGATCGA	ATCGATGACTAG
AATCGATGCTAT	TAGCACATGCGCT	ATCTTACAGTAA
ATACAGCTTCTAT	ATAGCCCTGAGT	CGTTAGATGATA
TAGATCGATGAA	CGTGTATGATAT	GCACATCGCTAT

Generate several thousand sequence reads per clone

Assemble

Green (2001)



February, 2001 Draft Sequence

The cover of Nature magazine from February 1, 2001. The title "nature" is in a large, white, serif font. Below it, "the human genome" is written in a smaller, white, sans-serif font. The background is a dark, textured image of a DNA double helix. Text on the cover includes "1 February 2001", "www.nature.com", "Nuclear fusion: Five-dimensional energy landscapes", "Seafloor spreading: The view from under the Arctic ice pack", "Cancer prospects: Sequence creates new opportunities", and "naturejobs genomics special".

The cover of Science magazine from February 16, 2001. The title "Science" is in a large, white, serif font. Below it, "THE HUMAN GENOME" is written in a smaller, white, sans-serif font. The background features a collage of human faces of various ethnicities and ages. Text on the cover includes "16 February 2001", "Vol. 291 No. 5507", "Pages 1145-1434 59", and "AMERICAN ASSOCIATION FOR THE ADVANCEMENT OF SCIENCE".

International Human Genome Sequencing Consortium (2001) **Venter et al. (2001)**

April, 2003 Completion



October, 2004 Publication

21 October 2004

International journal of science

nature

www.nature.com

Tetraodon to human
Evolutionary history in genome sequences

General relativity
Did the orbit move for you?

The human genome
Going the last mile

Antibiotics crisis
Market forces fail to deliver

Medical ethics
Choosing deafness

naturejobs think Finland

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

The sequence of the human genome encodes the genetic instructions for human physiology, as well as rich information about human evolution. In 2001, the International Human Genome Sequencing Consortium reported a draft sequence of the euchromatic portion of the human genome. Since then, the international collaboration has worked to convert this draft into a genome sequence with high accuracy and nearly complete coverage. Here, we report the final draft of this finishing process. The current genome sequence (build 30) contains 2.95 billion nucleotides interrogated by only 3.41 gaps. It covers ~99% of the euchromatic genome and is accurate to an error rate of ~1 error per 100,000 bases. Half of the remaining euchromatic gaps are associated with repetitive elements and will require focused work with new methods. The near-complete sequence, the first for a vertebrate, greatly improves the precision of biological analyses of the human genome including studies of gene number, birth and death. Finally, the human genome serves to encode only 20,000–25,000 protein-coding genes. The genome sequence reported here should serve as a firm foundation for biomedical research in the decades ahead.

The Human Genome Project (HGP) was launched in 1990 with the goal of obtaining a highly accurate sequence of the full majority of the euchromatic portion of the human genome. The initial work followed a two-pronged approach: (1) the mapping of the human and mouse genomes¹ to facilitate the study of inherited disease and provide a scaffold for genome assembly; and (2) the sequencing of regions with simple, repetitive genomes^{2,3} to serve as a method for method development and assist in interpreting the human genome. With success along both paths, the sequencing of the human genome itself eventually became feasible. The International Human Genome Sequencing Consortium (IHGSC), an open collaboration involving twenty centers in its constitution, was formed to carry out the component of the HGP.

In February 2001, the IHGSC⁴ and Celera Genomics⁵ each reported draft sequences providing a first overall view of the human genome. These sequences allowed systematic study of the human genome itself, including identification of genes, coordinate-based architecture of proteins, regional differences in genome composition, distribution and history of transposable elements, distribution of polymorphisms and relationships between genetic, molecular and physical distance. However, systematic knowledge of the human genome has enabled new tools and approaches that have enabled accelerated biomedical research.

Draft-level sequences, however, had important shortcomings. The IHGSC sequence, for example, covered ~90% of the euchromatic genome; it was arranged by ~100,000 gaps and the order and orientation of many sequences within local regions had not been established. The Celera sequence, while more complete in coverage of the euchromatic genome, operationally a finished draft, contained a high density of errors, including a high density of gaps being those reflective of all available technologies⁶ (see <http://www.genome.gov/2002/02/020202a.html>).

The goal of this finishing process was to convert the human genome to a higher quality sequence with improved accuracy and high-resolution information about the finishing process and finishing standards can be found in the Supplementary Information (Table 1) and <http://www.genome.gov/2002/02/020202b.html>.

In total, we generated a draft sequence from 39,268 large-insert clones (total length ~1.14 gigabases (Gb)) and finished the sequence from 43,742 of these clones (total length ~3.47 Gb). The clones contained primarily of bacterial artificial chromosomes

articles

across the goal of a complete human sequence. The number of gaps has been reduced 40-fold to only 341, most of which are associated with repetitive elements and will require new methods for resolution. The assembled near-complete genome sequence has an error rate of only ~1 error per 100,000 bases. It contains 2.95 billion nucleotides and covers ~99% of the euchromatic genome. This paper describes the current genome sequence and the process used to produce it, examines the accuracy and completeness of the sequence, and discusses biological evidence made possible by the sequence. We do not attempt here a comprehensive analysis of the content of the human genome, an initial analysis was previously reported⁴ and a series of papers is being written describing the individual chromosomes^{7,8}, including construction of genes and other features.

Current genome sequence

Human genome

The process of converting the initial draft sequence into a near-complete sequence is referred to as 'finishing'. It is a complex iterative process that generates simultaneously, at multiple scales, ranging from single nucleotides to the mapping of whole chromosomes. The fundamental challenge in this genomic region that are not well represented or easily resolved through random shotgun sequencing tend to be highly resolved in problematic regions. Finishing such regions required the development of special approaches, which varied substantially over time and varied among centers.

Finally, the finishing process involved two distinct components: (1) producing finished maps, consisting of contigs and overlap paths of overlapping large-insert clones spanning the euchromatic portion of each chromosome arm; and (2) producing finished clones, consisting of complete and accurate nucleotide sequences across each large insert clone. In practice, these two components were tightly intertwined in that progress in each often depended on results from the other. The components are described in Boxes 1 and 2.

Further information about the finishing process and finishing standards can be found in the Supplementary Information (Table 1) and <http://www.genome.gov/2002/02/020202b.html>.

International Human Genome Sequencing Consortium (2004)

CNN's #1 Medical Story of Past 25 Years

CNN.com

PRINT THIS

Powered by @clickability

Click to Print

SAVE THIS | EMAIL THIS | Close

Top 25: Medical stories

Human genome mapping ranks No. 1 in health news

Tuesday, March 29, 2005 Posted: 4:24 PM EST (2124 GMT)

(CNN) -- Much of the marvel of medicine has to do with discovery. Mapping the human genome, the complete sequence of DNA, gave scientists a blueprint for building a person, making it the No. 1 medical story, according to a distinguished panel CNN gathered to rank the top 25 medical stories of the past quarter-century.

Two men from two separate groups -- Francis Collins of the National Institutes of Health and Craig Venter of Celera Genomics Inc., a pharmaceutical-development company -- worked independently to discover the sequence of the human genome and identify the genes that it contains. This

April, 1953 → April, 2003

No. 4356 April 25, 1953 NATURE

MOLECULAR STRUCTURE OF NUCLEIC ACIDS

A Structure for Deoxyribose Nucleic Acid



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.



**All of the original goals of the
Human Genome Project have
been accomplished!**

What's Next?



A vision for the future of genomics research

A blueprint for the genomic era.

Francis S. Collins, Eric D. Green, Alan E. Guttmacher and Mark S. Guyer on behalf of the US National Human Genome Research Institute

The completion of a high-quality, comprehensive sequence of the human genome, in this effort, is an extraordinary year of the discovery of the double-helical structure of DNA, is a landmark event. The genome era is now a reality.

In contemplating a vision for the future of genomics research, it is appropriate to consider the remarkable progress that has brought us here. The milestone (Figure 1) shows a timeline of landmark accomplishments in genetics and genomics, beginning with Gregor Mendel's discovery of the laws of heredity and their rediscovery in the early days of the twentieth century. The sequencing of DNA and the hereditary material, determination of its structure, elucidation of the genetic code, development of recombinant DNA techniques, and establishment of increasingly automatable methods for DNA sequencing, all the way to the Human Genome Project (HGP) to begin in 1990 (see also www.nature.com/nature/ENAG0). Thanks to the vision of the original planners, and the creativity and determination of a legion of talented scientists who decided to make this project their overriding focus, all of the initial objectives of the HGP have now been achieved at least two years ahead of expectation, and a revolution in biology is now being begun.

The project's new research strategies and experimental techniques have generated a steady stream of ever larger and more complete genomic data sets that have been posted into public databases and have transformed the study of virtually all life processes. The genomic approach of technology development and large-scale generation of community resource data sets has introduced an important new dimension into biological and biomedical research. Intersectoral advances in genetics, comparative genomics, high-throughput biochemistry and microarray technology are providing biologists with a methodically improved repertoire of research tools that will allow the harnessing of organisms to health and disease to be analyzed and comprehended at an unprecedented level of molecular detail. Genome sequences, the bounded sets of information that guide biological development and function, lie at the heart of this revolution. In short, genomics has become a central and cohesive discipline of biological research.

The practical consequences of the emergence of this new field are widely apparent. Identification of the genes responsible for human mendelian disorders, once a herculean task requiring large research teams, many years of hard work, and an enormous outlay, can now be routinely accomplished

feature

in a few weeks by a single graduate student with access to DNA samples and associated phenotypes, an Internet connection to the public genome databases, a thermal cycler and a DNA sequencing machine. With the recent publication of a draft sequence of the mouse genome¹, identification of the mouse genome², identification of the mouse genome³, identification of the mouse genome⁴, a number of interesting mouse phenotypes has already been greatly simplified. Comparison of the human and mouse sequences shows that the proportion of the mammalian genome under evolutionary selection is more than twice that previously assumed.

Our ability to explore genome function is increasing in specificity as each subsequent genome is sequenced. Microarray technologies have catapulted many laboratory scientists studying the expression of one or two genes in a month to studying the expression of tens of thousands of genes in a single afternoon⁵. Critical opportunities for gene-based pre-emptive medicine, prediction of illness and adverse drug response are emerging at a rapid pace, and the therapeutic promise of genomics has taken on an exciting phase of expansion and exploration in the commercial sector⁶. The investment of the HGP in making the ethical, legal and social implications of these scientific advances has created a talented cohort of scholars in ethics, law, social science, clinical research, theology and public policy, and has already resulted in substantial increases in public awareness and the introduction of significant (but still incomplete) protections against misuse such as genetic discrimination (see www.genome.gov/ELSI/ethics).

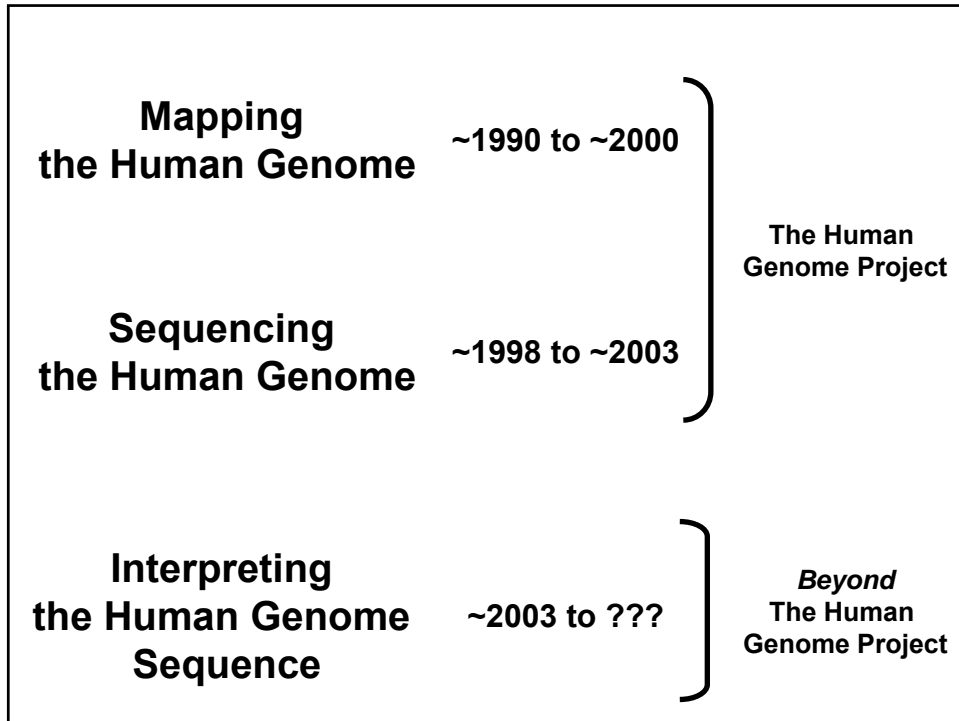
These accomplishments fulfill the expectations articulated in the 1990 report of the National Research Council, Mapping and Sequencing the Human Genome⁷. The successful completion of the HGP thus represents an opportunity to look forward and offer a blueprint for the future of genomics research over the next several years.

The vision presented here addresses a different world from that reflected in earlier plans published in 1990, 1993 and 1995 (refs. 15–17). These documents addressed the goals of the 1980 report, defining detailed paths towards the development of genome-

Collins et al. (2003)

Collins et al. (2003)

Page 26



Foundational Milestones in Genetics & Genomics



Darwin

1859



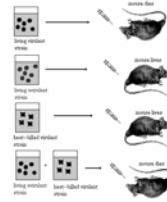
Mendel

1865



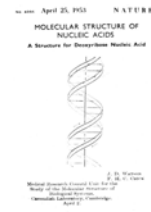
Miescher

1871



Avery

1944



Watson & Crick

1953

Comparing Genomes is Like Cryptography

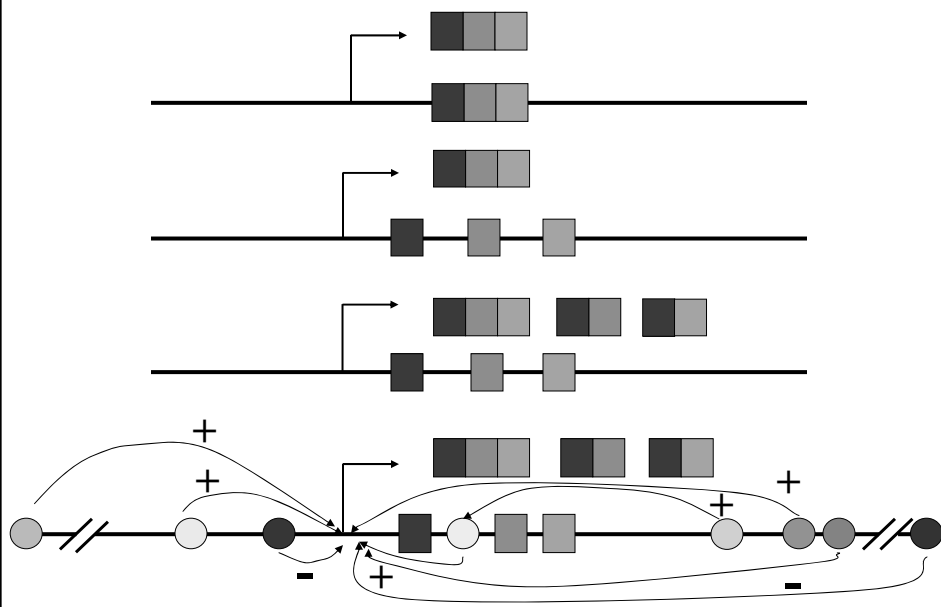
CKQEBHEREYTWASULSCZMEISDFOGETHEBLPGOODFQSTLKS TUFFRTAC

DLUCEHEREZBRTTOISAWNCDARJPTHERROFGOODERGHCLS TUFFBRHA

Functional Elements: Coding vs. Non-Coding

- **Coding Sequences (i.e., Genes)**
 - Relatively EASY to Identify
 - Mostly Know What to Look For
 - Complementary Data Sets Available (ESTs, cDNAs)
 - Ever-Improving Computational Gene Predictions
- **Non-Coding Functional Sequences**
 - HARD to Identify
 - Very Little Known About What to Look For
 - Virtually No Complementary Data Sets Available
 - Poor Computational Predictions

The Language of the Genome



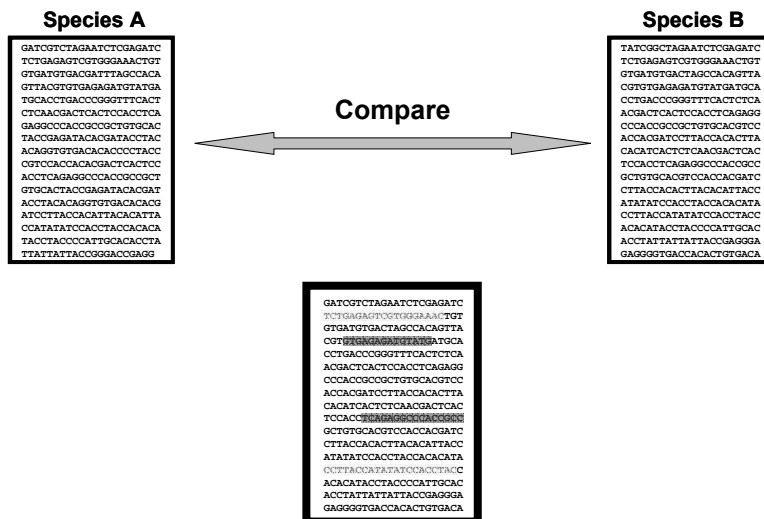
Functional Elements: Coding vs. Non-Coding

- **Coding Sequences (i.e., Genes)**
 - Relatively EASY to Identify
 - Mostly Know What to Look For
 - Complementary Data Sets Available (ESTs, cDNAs)
 - Ever-Improving Computational Gene Predictions
- **Non-Coding Functional Sequences**
 - HARD to Identify
 - Very Little Known About What to Look For
 - Virtually No Complementary Data Sets Available
 - Poor Computational Predictions


Major role for comparative sequence analysis will be the identification of functionally important, non-coding sequences

Comparative Sequence Analysis


Using the Experiments of Evolution to Decode the Human Genome




Vertebrate Genome Sequences




Mouse




Rat




Chicken




Chimpanzee




Dog




Macaque



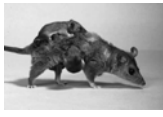
Orangutan




Marmoset




Cow



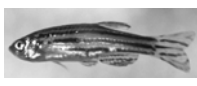
Monodelphis



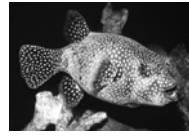
Platypus



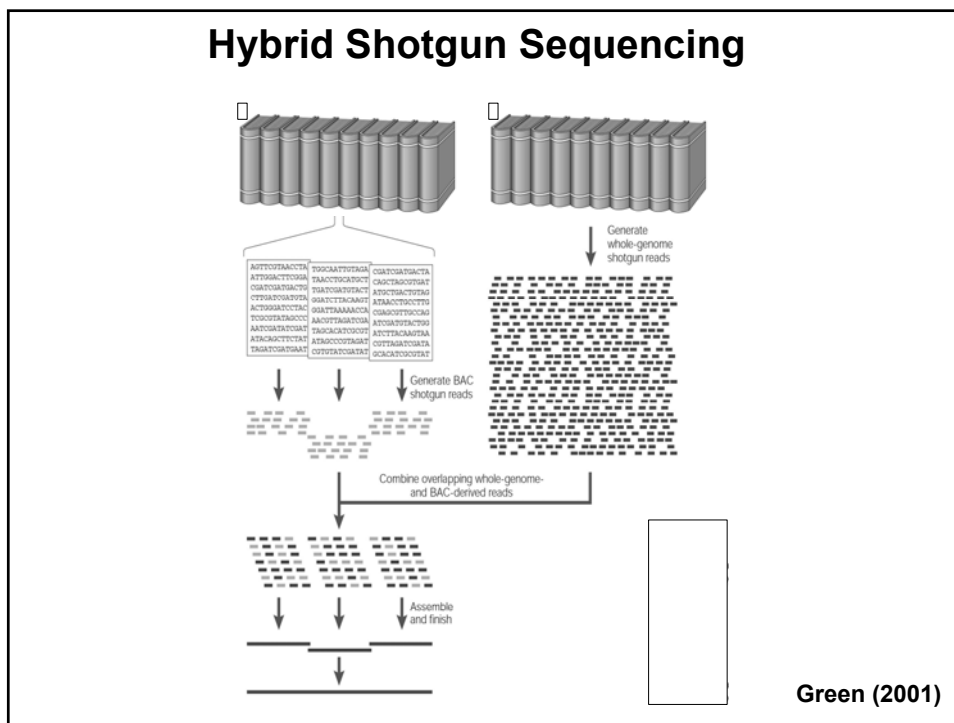
Xenopus



Zebrafish



Pufferfish

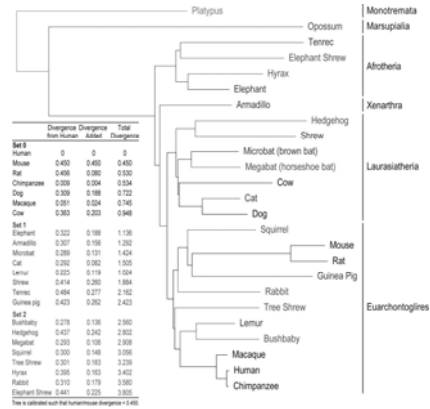


Low-Redundancy, Whole-Genome Shotgun Sequencing

An initial strategy for the systematic identification of functional elements in the human genome by low-redundancy comparative sequencing

Elliott H. Margulies^{1*}, Jade Vinson^{1*}, NISC Comparative Sequencing Program^{1,2*}, Webb Miller¹, David B. Jaffe¹, Kerstin Lindblad-Toh¹, Jean Chang¹, Eric D. Green^{1,3}, Eric S. Lander¹, James C. Mullikin^{1,3*}, and Michele Clamp^{1,3*}

Margulies EM et al. (2005)



Landscape of Vertebrate Genome Sequencing

Human	=====
Mouse	=====
Rat	=====
Pufferfish	=====
Zebrafish	=====
Chicken	=====
Chimpanzee	=====
Dog	=====
Cow	=====
Xenopus	=====
Monodelphis	=====
Macaque	=====
Platypus	=====
Marmoset	=====
Orangutan	=====
Armadillo
Elephant
Tenrec
Rabbit
Cat
Shrew
Guinea Pig
Hedgehog
(and others...)	

Multi-Species Sequence Comparisons

<pre>GATGCTCTAGAACTCTCG AGATCTCTGAGAGTCGT GGGAAACTGTGTGATGT GACTAGCCACGTTACG TGTGAGAGATGTATGAT GCACCTGACCCGGGTTT CAGCTTCAAGACTCAC TCCACCTCAGAGGCCCA CCGCCCCTGTGCAGCTC CACACAGACCTTACCA CAGTTACAGATTACCAT ATATCCACCTACCACAC ATACTACCCCATTTGCA CAGCTATATATATACG</pre>	<pre>GATGCTCTAGAACTCTCG AGATCTCTGAGAGTCGT GGGAAACTGTGTGATGT GACTAGCCACGTTACG TGTGAGAGATGTATGAT GCACCTGACCCGGGTTT CAGCTTCAAGACTCAC TCCACCTCAGAGGCCCA CCGCCCCTGTGCAGCTC CACACAGACCTTACCA CAGTTACAGATTACCAT ATATCCACCTACCACAC ATACTACCCCATTTGCA CAGCTATATATATACG</pre>	<pre>GATGCTCTAGAACTCTCG AGATCTCTGAGAGTCGT GGGAAACTGTGTGATGT GACTAGCCACGTTACG TGTGAGAGATGTATGAT GCACCTGACCCGGGTTT CAGCTTCAAGACTCAC TCCACCTCAGAGGCCCA CCGCCCCTGTGCAGCTC CACACAGACCTTACCA CAGTTACAGATTACCAT ATATCCACCTACCACAC ATACTACCCCATTTGCA CAGCTATATATATACG</pre>	<pre>GATGCTCTAGAACTCTCG AGATCTCTGAGAGTCGT GGGAAACTGTGTGATGT GACTAGCCACGTTACG TGTGAGAGATGTATGAT GCACCTGACCCGGGTTT CAGCTTCAAGACTCAC TCCACCTCAGAGGCCCA CCGCCCCTGTGCAGCTC CACACAGACCTTACCA CAGTTACAGATTACCAT ATATCCACCTACCACAC ATACTACCCCATTTGCA CAGCTATATATATACG</pre>	<pre>GATGCTCTAGAACTCTCG AGATCTCTGAGAGTCGT GGGAAACTGTGTGATGT GACTAGCCACGTTACG TGTGAGAGATGTATGAT GCACCTGACCCGGGTTT CAGCTTCAAGACTCAC TCCACCTCAGAGGCCCA CCGCCCCTGTGCAGCTC CACACAGACCTTACCA CAGTTACAGATTACCAT ATATCCACCTACCACAC ATACTACCCCATTTGCA CAGCTATATATATACG</pre>	<pre>GATGCTCTAGAACTCTCG AGATCTCTGAGAGTCGT GGGAAACTGTGTGATGT GACTAGCCACGTTACG TGTGAGAGATGTATGAT GCACCTGACCCGGGTTT CAGCTTCAAGACTCAC TCCACCTCAGAGGCCCA CCGCCCCTGTGCAGCTC CACACAGACCTTACCA CAGTTACAGATTACCAT ATATCCACCTACCACAC ATACTACCCCATTTGCA CAGCTATATATATACG</pre>	<pre>GATGCTCTAGAACTCTCG AGATCTCTGAGAGTCGT GGGAAACTGTGTGATGT GACTAGCCACGTTACG TGTGAGAGATGTATGAT GCACCTGACCCGGGTTT CAGCTTCAAGACTCAC TCCACCTCAGAGGCCCA CCGCCCCTGTGCAGCTC CACACAGACCTTACCA CAGTTACAGATTACCAT ATATCCACCTACCACAC ATACTACCCCATTTGCA CAGCTATATATATACG</pre>
---	---	---	---	---	---	---

HUMAN

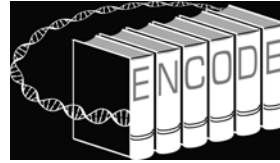
Multi-Species Conserved Sequences (MCSs)

Margulies et al. (2003)

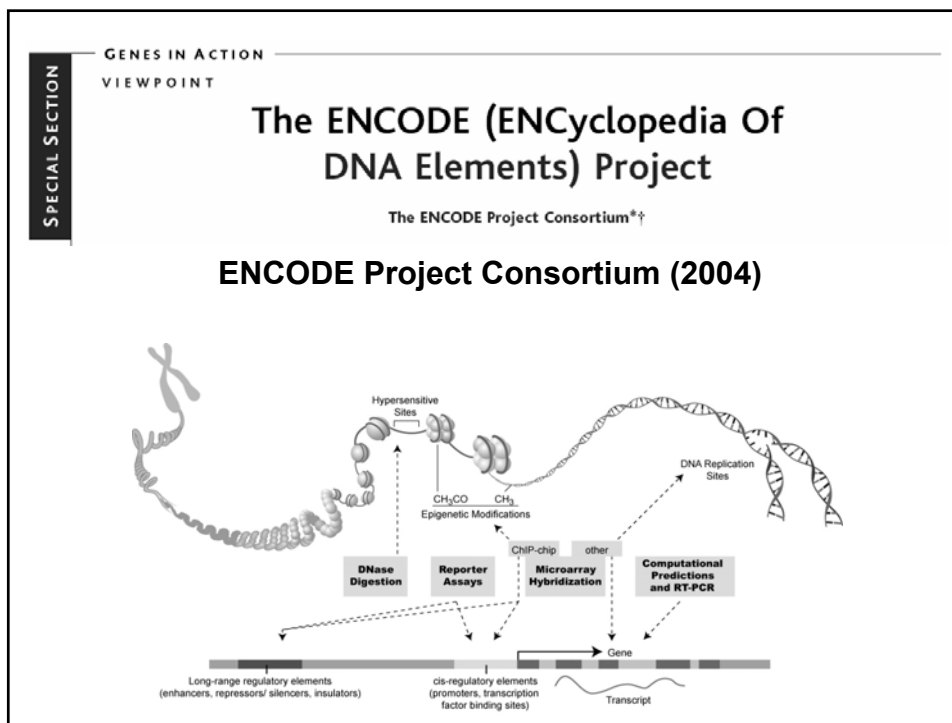
Future Genomes to Sequence???

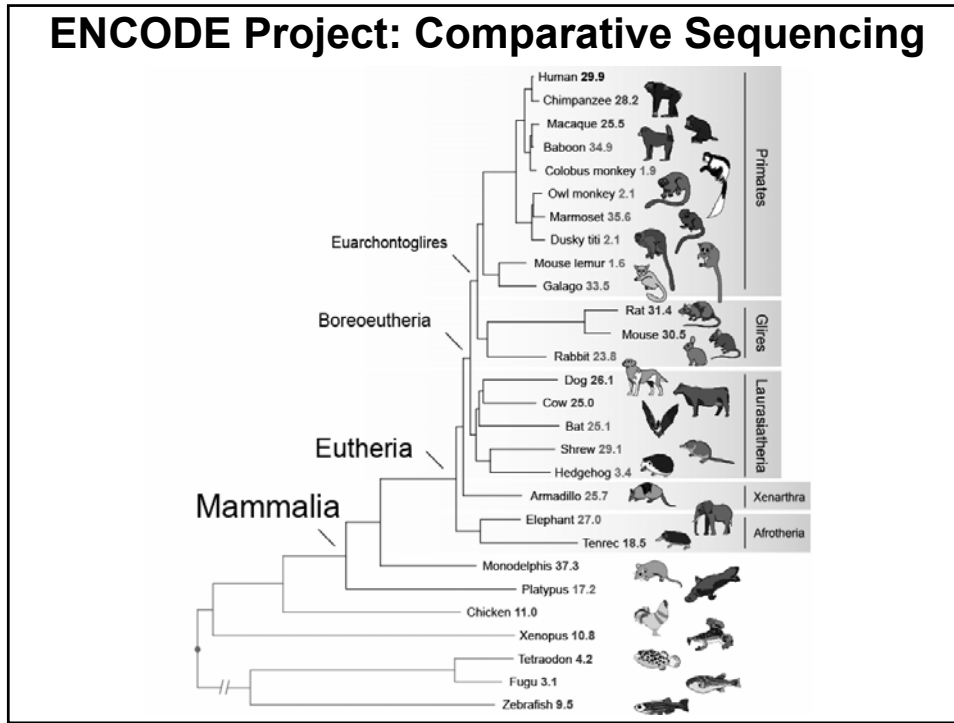


ENCODE Project



- ENCODE: ENCyclopedia Of DNA Elements
- Goal: Compile a *Comprehensive Encyclopedia* of All Functional Elements in the Human Genome
- Initial Pilot Project: 1% of Human Genome
- Apply Multiple, Diverse Approaches to Study and Analyze that 1% in a Consortium Fashion





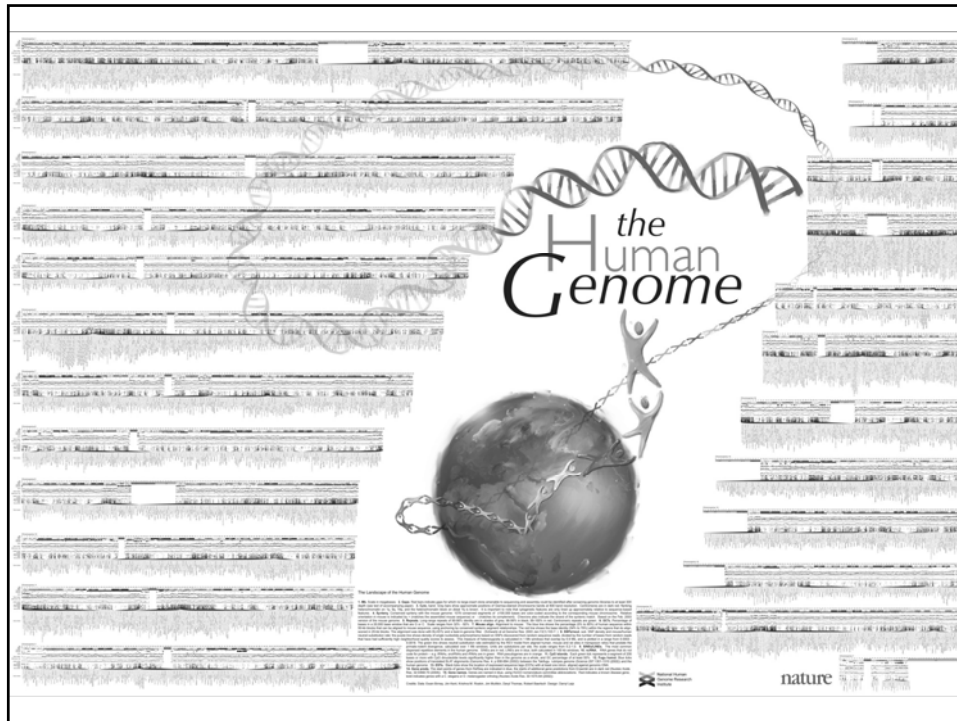
ENCODE Project: Web Sites

The screenshot shows the ENCODE Project website. The main heading is "The ENCODE Project: ENCyclopedia of DNA Elements". Below it, there are several sections: "ENCyclopedia of DNA Elements", "ENCyclopedia of Target Regions (January 2004)", and "ENCyclopedia of Conservation". The "ENCyclopedia of Target Regions" section includes a table of "Stratified Random Picks" with columns for "Region", "Non-Exonic Conservation (%)", "GC Density (%)", and "GC Density (%)".

The screenshot shows the UCSC Genome Browser interface. The main heading is "UCSC Genome Browser on Human July 2003 Freeze". The browser displays a genomic track for a region on chromosome 1. The track includes various annotations such as "ENCODE Regions", "ENCODE Regions", "ENCODE Regions", and "ENCODE Regions". The browser also shows a list of "ENCODE Regions" with columns for "Region", "Description", "Chr", and "Size".

genome.gov/ENCODE

genome.ucsc.edu/ENCODE



Human Genome Sequence

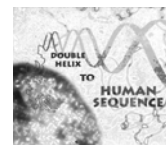
>\$1,000,000,000



~\$100,000



~\$1,000



WIRED TOOLS 2K3 [TOYS] by Chris Baker

MY FIRST DNA SEQUENCE

The same tech requirements that are reshaping green-up gadgets are revolutionizing kids' toys. Now, thanks to a simple, user-friendly, remote control truck, kids with programmable robots, and great ear on the latest electronic gizmos, are excited enough for your grandkids to take the Discovery Kids DNA Explorer home. Junior scientists extract and read real deoxyribonucleic acid. As third grade science projects go, this is light years beyond the of baking soda volcano. Next stop: cloning fish.

[DNA Explorer Ages 10 and up. \$99. www.discovery.com]



Prep the Specimen
Series extracting DNA, your young Dr. Frankenstein has to pick a specimen and prep it. The experiment works on all kinds of food, like corn, beans, or even (gross!) chicken feet. The kit includes freeze-dried ground peas (that DNA is easy to extract), isopropyl alcohol, and salt. When all three are put in a beaker and mixed with distilled water, the giant, cellular structure starts to break down.

Separate the Oils
After transferring the mixture to a test tube, Dr. Frankenstein needs to add dish soap. The tube goes inside a splitter-proof magnetic mixer and centrifuge, which pulls up the oils and separates the DNA from the saline liquid. After 10 seconds, Dr. Frankenstein is a pinch of anionics, with alcohol, and the DNA strands float to the surface, where they can be harvested with the "DNA Hook."

Zap the Molecules
To read the code of DNA, Dr. Frankenstein needs to zap up a bit of conductive gel. It reads from 100 buffer and separate possible. The gel goes into a battery-operated electrophoresis chamber, where it's zapped with a tool to make ditches for the fragmented genes. The molecules are transferred with a pipette. A zap of electricity sends the molecules - which are negatively charged - moving through the gel.

Unravel the Mystery
A couple of hours later, Dr. Frankenstein will see a few strips of stain (see second) - a basic reveal (not to confuse the gene) - or any other - genetic blueprint. Which year breakthrough? Frankenstein may figure out who Daddy really is (or isn't).

188-120003-0000

Discovery CHANNEL Enter Email **Gift Cards**

STORE DVD, VHS & BOOKS TELESCOPES & BINOCULARS TOYS & GAMES HOME & OFFICE ELECTRONICS & GADGETS FAN GEAR & CLOTHING

Carting Quick Order Save Items Gift Finder Store Locator NEW SEARCH

Discovery Pro Cycling Team Jersey, T-Shirts, Framed Prints and More... SHOP NOW

Home > Search Results > Discovery DNA Explorer Kit

EXCLUSIVE
Discovery DNA Explorer Kit
#698795

CUSTOMER RATING: ★★★★★

\$79.95
IN STOCK

Include all required batteries **\$20.00**

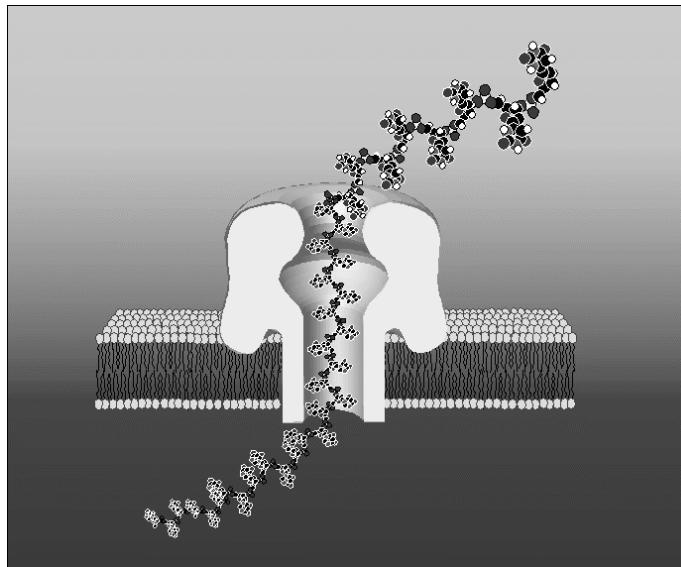
ADD TO CART

Shipping Information: Kit is available for shipment to the U.S. & Canada.

Product Detail:
Winner of the 2004 Parents' Choice Gold Award, Learning Magazine's Teachers' Choice Award and named to Popular Science's 2003 Best of What's New List.

Please Note: Two of the 10 experiments included in the Kit require the

En Route to the \$1000 Genome



Genome sequencing in microfabricated high-density picolitre reactors

Marcel Margulies¹, Michael Egholm¹, William E. Altman¹, Said Attiya¹, Joel S. Bader¹, Lisa A. Bemben¹, Jan Berka¹, Michael S. Braverman¹, Yi-Ju Chen¹, Zhoutao Chen¹, Scott B. Dewell¹, Lei Du¹, Joseph M. Fierro¹, Xavier V. Gomes¹, Brian C. Godwin¹, Wen He¹, Scott Helgeson¹, Chun He Ho¹, Gerard P. Irzyk¹, Szilveszter C. Jando¹, Maria L. I. Alenquer¹, Thomas P. Jarvie¹, Kshama B. Jirage¹, Jong-Bum Kim¹, James R. Knight¹, Janna R. Lanza¹, John H. Leamon¹, Steven M. LeKowitz¹, Ming Lei¹, Jing Li¹, Kenton L. Lohman¹, Hong Lu¹, Vinod B. Makhijani¹, Keith E. McDade¹, Michael P. McKenna¹, Eugene W. Myers¹, Elizabeth Nickerson¹, John R. Noble¹, Ramona Plant¹, Bernard P. Puc¹, Michael T. Ronan¹, George T. Roth¹, Gary J. Sarkis¹, Jan Fredrik Simons¹, John W. Simpson¹, Maitreyan Srinivasan¹, Karrie R. Tartaro¹, Alexander Tomasz², Karl A. Vogt¹, Greg A. Volkmer¹, Shally H. Wang¹, Yong Wang¹, Michael P. Weiner¹, Pengguang Yu¹, Richard F. Beigley¹ & Jonathan M. Rothberg¹

Margulies M et al. (2005)

Solexa Ltd

Simon Bennett PhD
Business Development Director

Solexa Ltd is developing an integrated system, based on a breakthrough single molecule sequencing technology, to address a US\$2 billion market that is expected to grow



Bennett (2004)

Toward the \$1000 human genome

Simon T Bennett,
Cedric Barnes,
Anthony Cox,
Lisa Chavez &
Clive Brena¹

Revolutionary new technologies, capable of transforming the economics of sequencing, are providing an unparalleled opportunity to analyze human genetic variation comprehensively at the whole-genome level within a realistic timeframe and at affordable costs. Current estimates suggest that it would cost somewhere in the region of US\$30 million to sequence



Bennett et al. (2005)

Accurate Multiplex Polony Sequencing of an Evolved Bacterial Genome

Jay Shendure^{1,8,†}, Gregory J. Porreca^{1,8,†}, Nikos B. Reppas¹, Xiaoxia Lin¹, John P. McCutcheon^{2,3}, Abraham M. Rosenbaum¹, Michael D. Wang¹, Kun Zhang¹, Robi D. Mitra², George M. Church¹

Shendure 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, forensics, and diagnostic and applied therapeutics. Several emerging technologies show promise of delivering next-generation solutions 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 to highlight recent advances and potential challenges these technologies face in their development for strand DNA sequencing.

More than just a mapping and sequencing endeavor, the Human Genome Project (HGP) has altered the mindset and approach to many basic and applied research efforts. Early disruption and controversy (Kirkland 1999; Lusa et al. 1998; Roberts 1998); Foa et al. 1996) were soon laid to rest by well-developed strategies (Roberts 1999a; Collins and Galis 1995; Collins et al. 1998) that led to the successful completion of mankind's largest biology project. As the cost of the HGP was including developments that advanced the pace of sequencing a mammalian genome from years to months. Along the way, numerous strategies emerged that hold promise for rapid, efficient, and accurate delivery of DNA sequence information. For the HGP, a break-through approach was adopted for completing the job by coupling the core technologies of Sanger sequencing and fluorescence detection. The completion of the sequencing phase could not have been accomplished without major innovations in microfluidic systems engineering, fluorescence dye development, capillary electrophoresis, automation, robotics, informatics, and process management. The result was completion of a high-quality, reference sequence of the human genome in April, 2003 (Collins et al. 2003), marking the 50-year anniversary of the discovery of the double-helix structure for many outside the genome community, that heroic milestone signalled the end of this international scientific project, but for the rest of us, it only marked the beginning of things to come.

The need for sequencing has never been greater than it is today, with applications spanning diverse research sectors including comparative genomics and evolution, forensics, epidemiology, and applied medicine for diagnostics and therapeutics. Arguably, the strongest impetus for emerging sequencing is the quest for identification and interpretation of human sequence variation in its relation to health and disease. The most common form of variation is the single nucleotide polymorphism (SNP), although two unrelated people share, on average, 99.9% sequence identity (i.e., one difference in a thousand base pairs), the average occurrence of an SNP in the general population is once every five hundred base pairs. As such, more than nine million unique SNPs have been cataloged in the public database, dbSNP (Conrad and Willerson 2005), with many more expected to be found in large-scale sequencing efforts.

A great deal of attention has been focused on common SNPs

with a minor allele frequency >5% and their potential role in common disease (Lander 1996; Rich and McKenzie 1996; Collins et al. 1997). Recent, large-scale genotyping efforts of these common SNPs have shown that much of the human genome can be parsed into common haplotype blocks (Daly et al. 2001; Fritsch et al. 2001; Gabriel et al. 2002). The International HapMap Consortium (2003) was formed to characterize common patterns of sequence variation by determining allele frequencies and the degree of association between SNPs among geographically distinct groups, leading to the identification of "tagSNPs" for genome-wide, disease-based association studies. With this method of characterization, however, rare SNPs/haplotypes may be overlooked, as highlighted by Lu et al. (2005), who described an association of rare "mutator/haplotypes" with congenital

A shift in large-scale strategies from genotyping to resequencing is currently taking place to explore the significance of less-common SNPs to human biology and disease. The "re" in this approach is the resequencing of individual genomes related to a reference genome for de novo SNP discovery and comparative genomics application. The ENCODE Project Consortium (2004) has devoted significant efforts toward resequencing megabase-sized blocks of the human genome. Consequently, genome centers are now offering at least 100x, 200x, or 300x coverage, which currently translates to ~8% capacity, to resequencing hundreds to thousands of gene regions. This increase in attention for high-throughput resequencing will greatly facilitate studies to determine the genetic basis of susceptibility to common disease, cancer biology, and disease association in model and nonmodel organisms.

Current sequencing technologies are too expensive, labor intensive, and time consuming for broad application to human sequence variation studies. Genome center costs (calculated on the basis of dollars per 1000 bp, bases defined below) can be generally divided into the categories of instrumentation, personnel, reagents and materials, and covered expenses. Currently, these costs are operating at less than one dollar per 1000 bp, bases, with at least 50% of the cost from DNA sequencing instrumentation alone. Developments in novel detection methods, miniaturization in instrumentation, microfluidic separation technologies, and an increase in the number of assays per run will soon likely have the biggest impact on reducing cost. It should be emphasized, however, that new sequencing strategies will be needed to rise above high-throughput platform efficiency. In September, 2004, the National Human Genome Res-

115176-1176 ©2005 by Cold Spring Harbor Laboratory Press, ISSN 1089-931X/05, www.genome.org

Genome Research 17:67 www.genome.org

Metzker (2005)

TSUNAMI SCIENCE: ONE YEAR AFTER THE WAVE THAT ROCKED THE WORLD

SCIENTIFIC AMERICAN

Alternatives to Toxic Tests on Animals

JANUARY 2006 \$4.99

WWW.SCIAM.COM

Know Your DNA

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

The Hazy Origin of Brown Dwarf Stars

Winning Tricks of the Racing Robots

Does Motherhood Make Women Smarter?



Church (2006)

Page 39

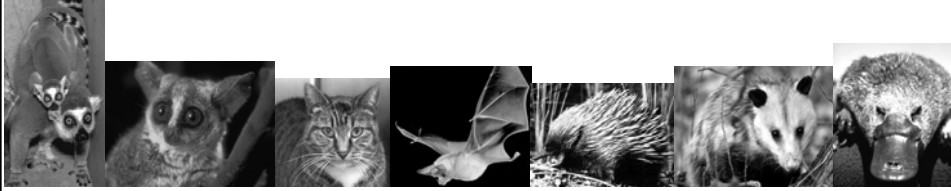
DNA Sequencing Technologies

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

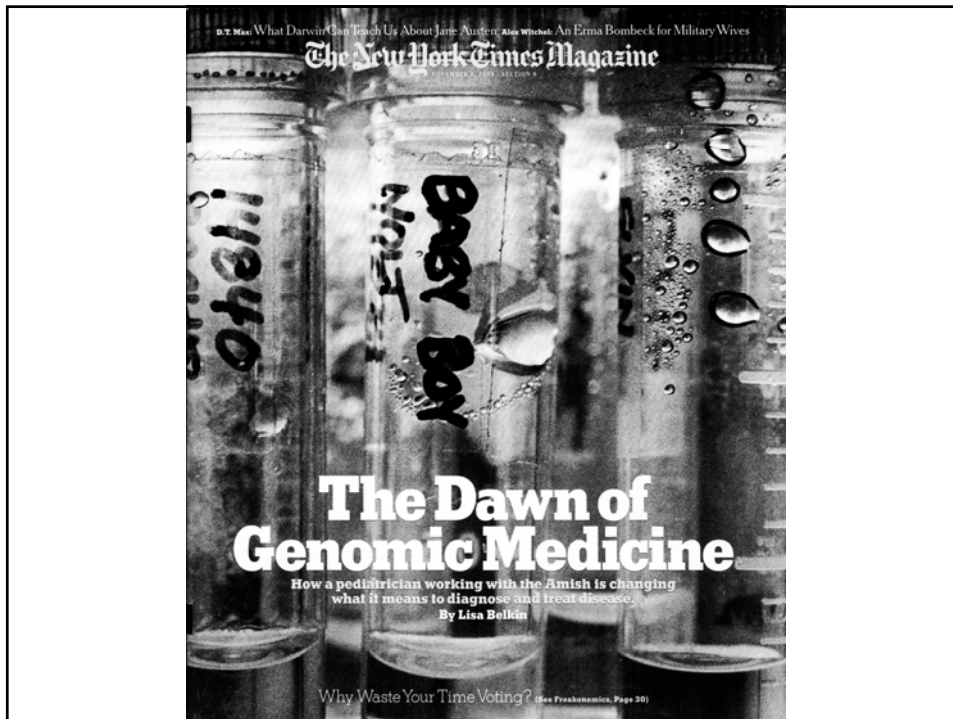
Realities of New DNA Sequencing Technologies...

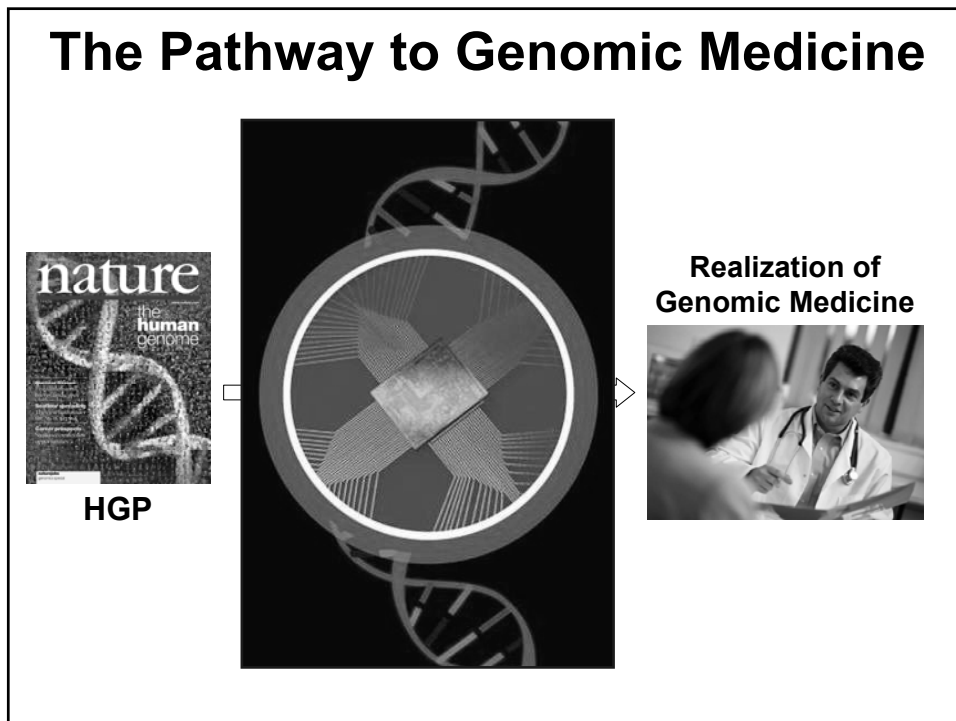
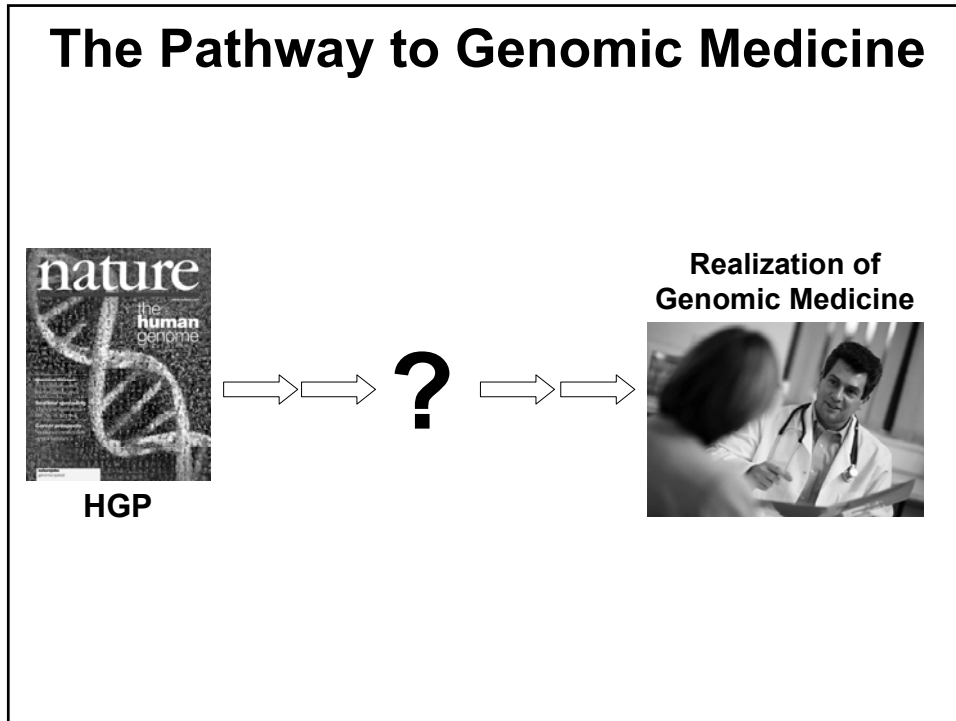


“Inter-Species” Comparisons



“Intra-Species” Comparisons





The Current Big Challenges...

- **Defining “Saturation Points” in Terms of Information Gained by Comparative Sequence Analyses**
- **Achieving the “\$1000 Genome”**
- **Large-Scale Deployment of Medical Sequencing**

The Human Genome Sequence to Genomic Medicine...



...from base pairs to bedside.

Bibliography

- Adams MD et al. (2000). The genome sequence of *Drosophila melanogaster*. *Science* 287:2185-2195.
- Bennett S (2004). Solexa Ltd. *Pharmacogenomics* 5:433-438.
- Bennett ST (2005). Toward the \$1000 human genome. *Pharmacogenomics* 6:373-382.
- Birren B et al. (1998). Bacterial artificial chromosomes. In *Genome Analysis: A Laboratory Manual, Vol. 3 Cloning systems* (B Birren et al., eds.; Cold Spring Harbor Laboratory Press), pp. 241-295.
- C. elegans* Sequencing Consortium (1998). Genome sequence of the nematode *C. elegans*: a platform for investigating biology. *Science* 282:2012-2018.
- Chimpanzee Sequencing and Analysis Consortium (2005). Initial sequence of the chimpanzee genome and comparison with the human genome. *Nature* 437:69-87.
- Church GM (2006). Genomes for all. *Sci Am* 294:46-54.
- Collins FS et al. (2003). A vision for the future of genomics research: a blueprint for the genomic era. *Nature* 422:835-847.
- ENCODE Project Consortium (2004). The ENCODE (ENCyclopedia Of DNA Elements) Project. *Science* 306:636-640.
- Gerhard DS et al. (2004). The status, quality, and expansion of the NIH full-length cDNA project: the Mammalian Gene Collection (MGC). *Genome Res* 14:2121-2127.
- Goffeau A et al. (1997). The Yeast Genome Directory. *Nature* 387S:1-105.
- Gordon D et al. (1998). Consed: a graphical tool for sequence finishing. *Genome Res* 8:195-202.
- Green ED (2001). Strategies for the systematic sequencing of complex genomes. *Nature Rev Genet* 2:573-583.
- Green ED et al. (1998). Yeast artificial chromosomes. In *Genome Analysis: A Laboratory Manual, Vol. 3 Cloning systems* (B Birren et al., eds.; Cold Spring Harbor Laboratory Press), pp. 297-565.
- Gregory SG et al. (1997). Genome mapping by fluorescent fingerprinting. *Genome Res* 7:1162-1168.
- Hillier LW et al. (2004). Sequence and comparative analysis of the chicken genome provide unique perspectives on vertebrate evolution. *Nature* 432:695-716.
- International Human Genome Sequencing Consortium (2001). Initial sequencing and analysis of the human genome. *Nature* 409:860-921.

- International Human Genome Sequencing Consortium (2004). Finishing the euchromatic sequence of the human genome. *Nature* 431:931-945.
- Lindblad-Toh K et al. (2005). Genome sequence, comparative analysis and haplotype structure of the domestic dog. *Nature* 438:803-819.
- Margulies EM et al. (2003). Identification and characterization of multi-species conserved sequences. *Genome Res* 13:2507-2518.
- Margulies EM et al. (2005). An initial strategy for the systematic identification of functional elements in the human genome by low-redundancy comparative sequencing. *Proc Natl Acad Sci* 102:4795-4800.
- Margulies M et al. (2005). Genome sequencing in microfabricated high-density picolitre reactors. *Nature* 437:376-380.
- Marra MA et al. (1997). High throughput fingerprint analysis of large-insert clones. *Genome Res* 7:1072-1084.
- Messing J and Llaca V (1998). Importance of anchor genomes for any plant genome project. *Proc Natl Acad Sci* 95:2017-2020.
- Metzker ML et al. (2005). Emerging technologies in DNA sequencing. *Genome Res* 15:1767-1776.
- Mouse Genome Sequencing Consortium (2002). Initial sequencing and comparative analysis of the mouse genome. *Nature* 420:520-562.
- Rat Genome Sequencing Project Consortium (2004). Genome sequence of the Brown Norway rat yields insights into mammalian evolution. *Nature* 428:493-521.
- Shendure J et al. (2005). Accurate multiplex polony sequencing of an evolved bacterial genome. *Science* 309:1728-1732.
- Thomas JW et al. (2003). Comparative analyses of multi-species sequences from targeted genomic regions. *Nature* 424:788-793.
- Venter JC et al. (2001). The sequence of the human genome. *Science* 291:1304-1351.
- Wilson RK and Mardis ER (1997). Fluorescence-based DNA sequencing. In *Genome Analysis: A Laboratory Manual, Vol. 1 Analyzing DNA* (B Birren et al., eds.; Cold Spring Harbor Laboratory Press), pp. 301-395.
- Wilson RK and Mardis ER (1997). Shotgun sequencing. In *Genome Analysis: A Laboratory Manual, Vol. 1 Analyzing DNA* (B Birren et al., eds.; Cold Spring Harbor Laboratory Press), pp. 397-454.