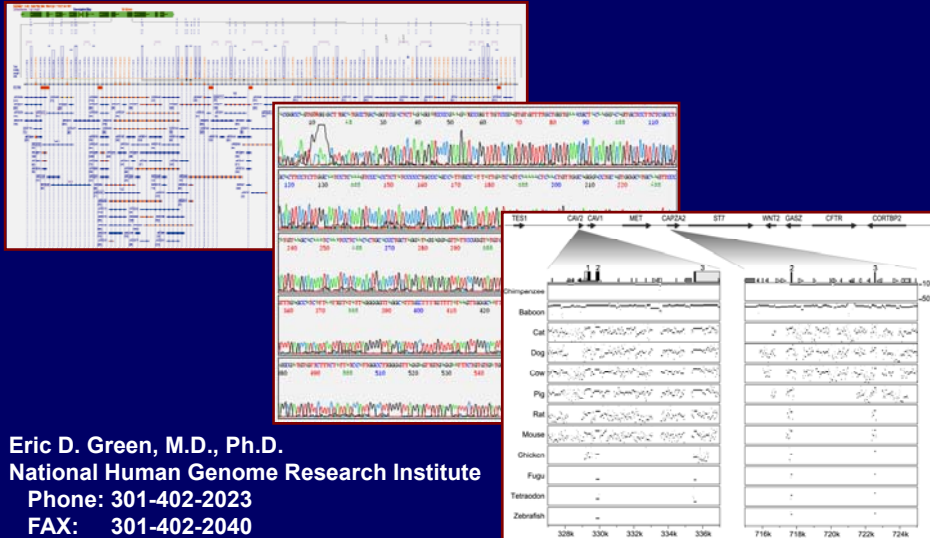


Techniques for Analyzing Genomes I



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



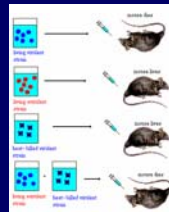
Mendel

1865



Miescher

1871



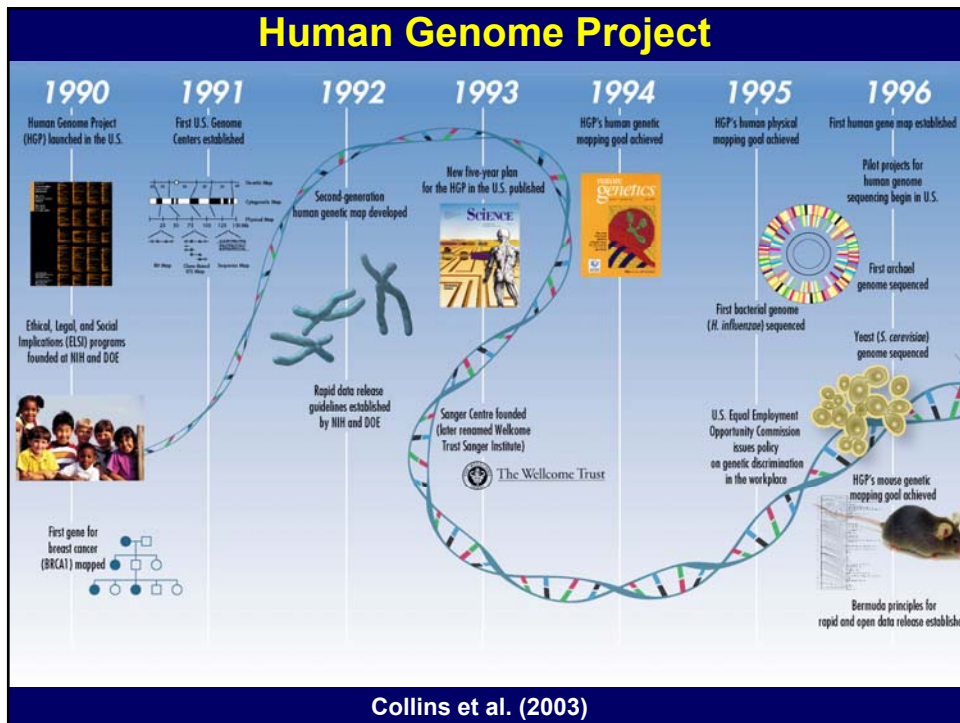
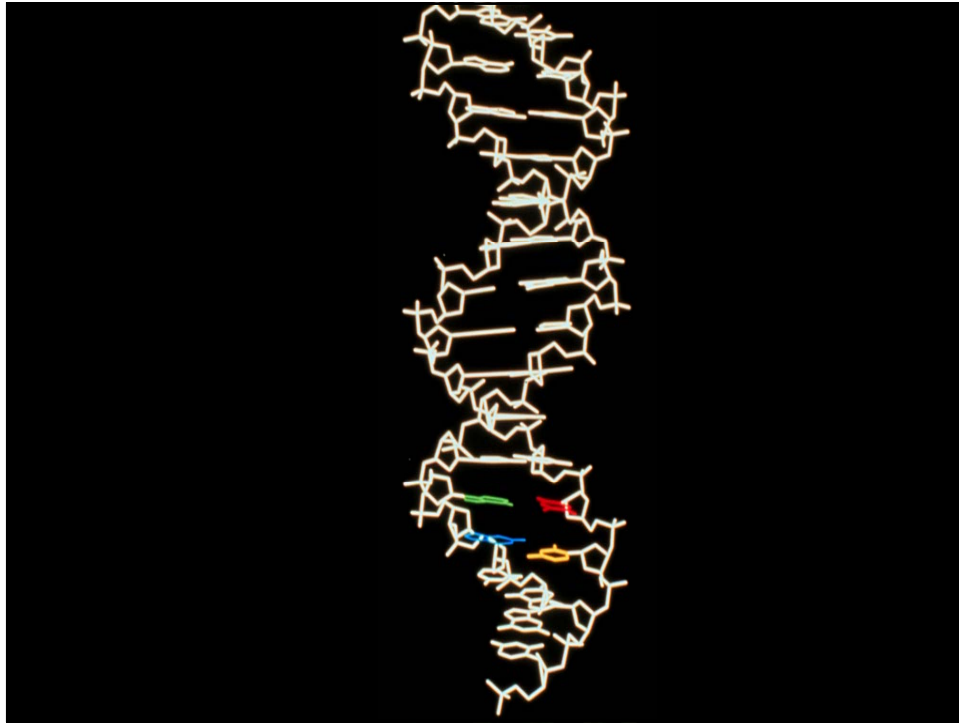
Avery

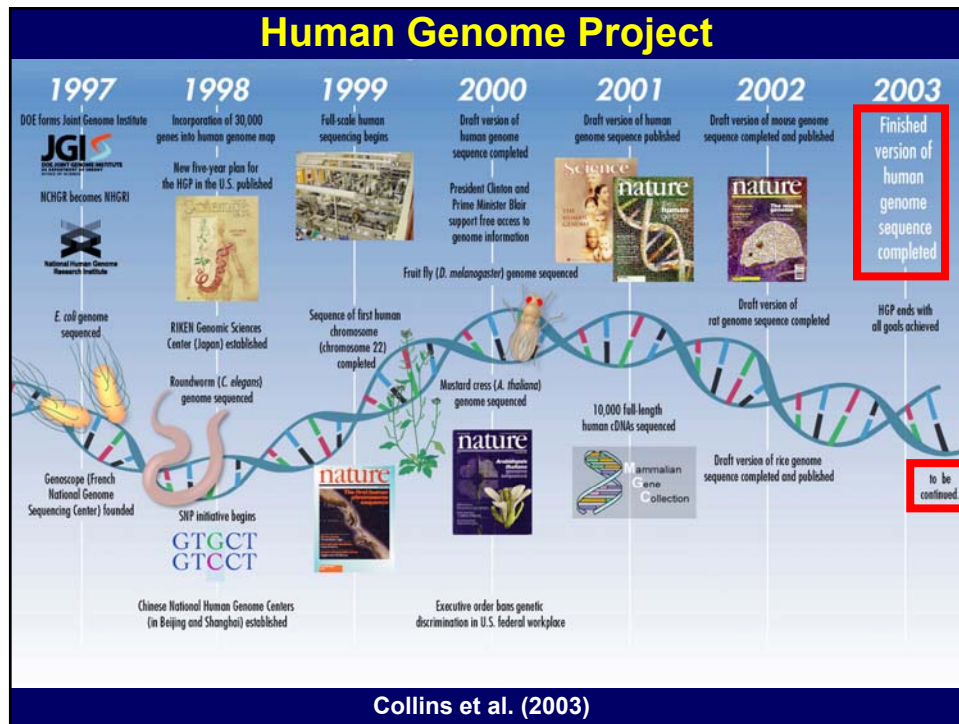
1944



Watson & Crick

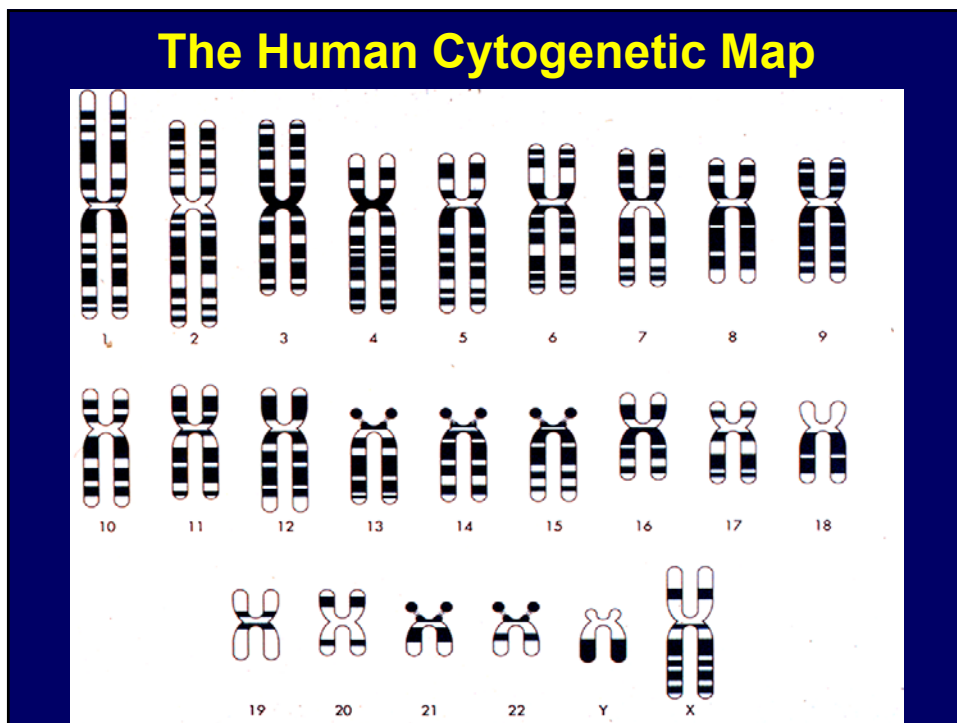
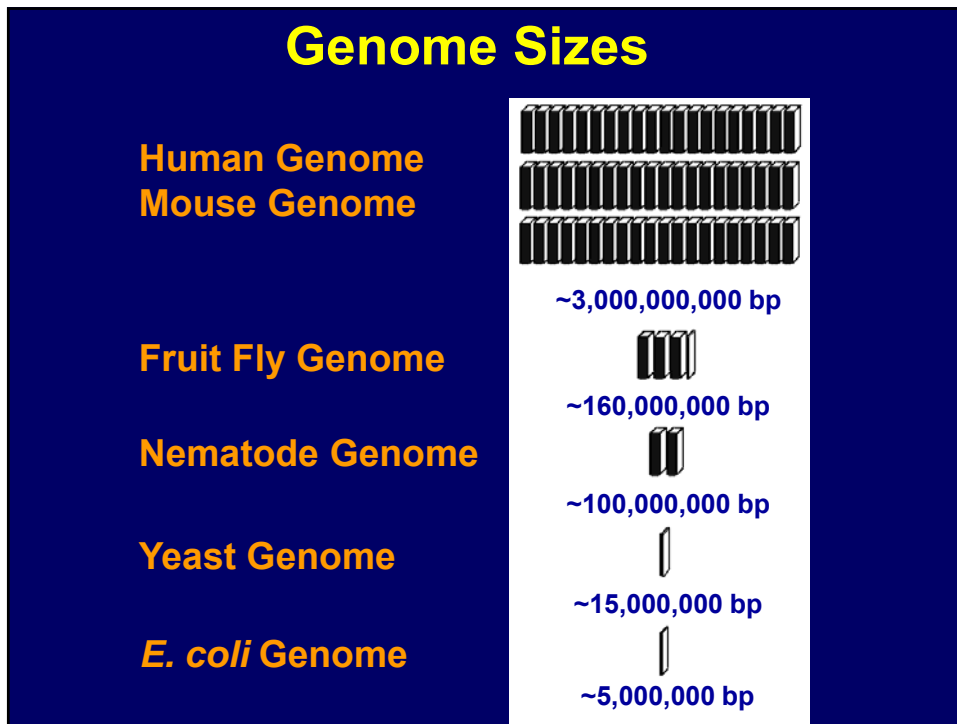
1953

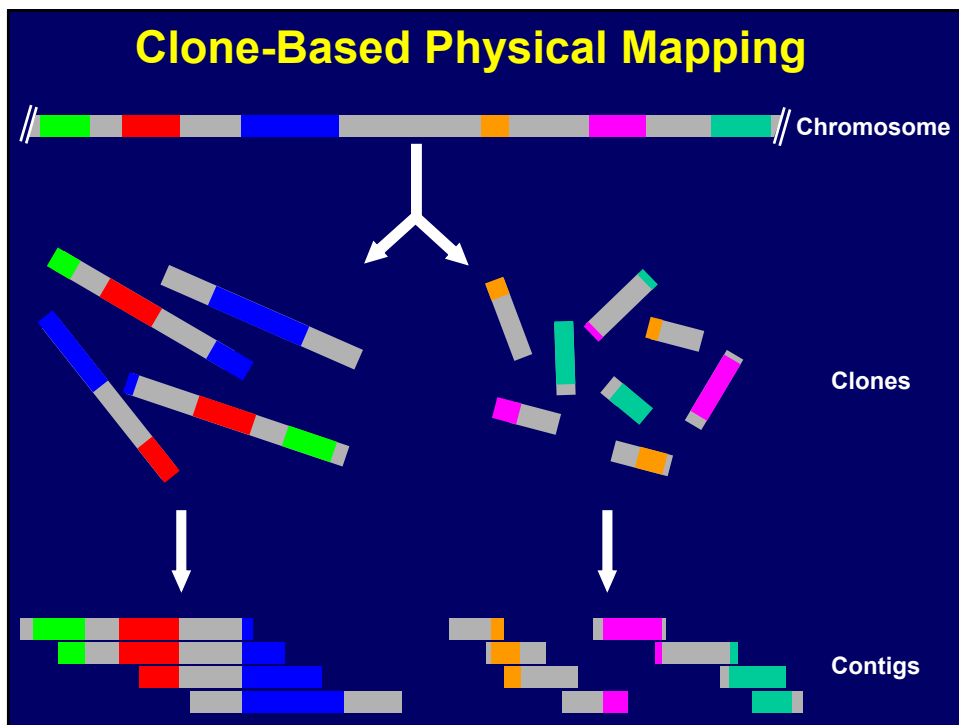
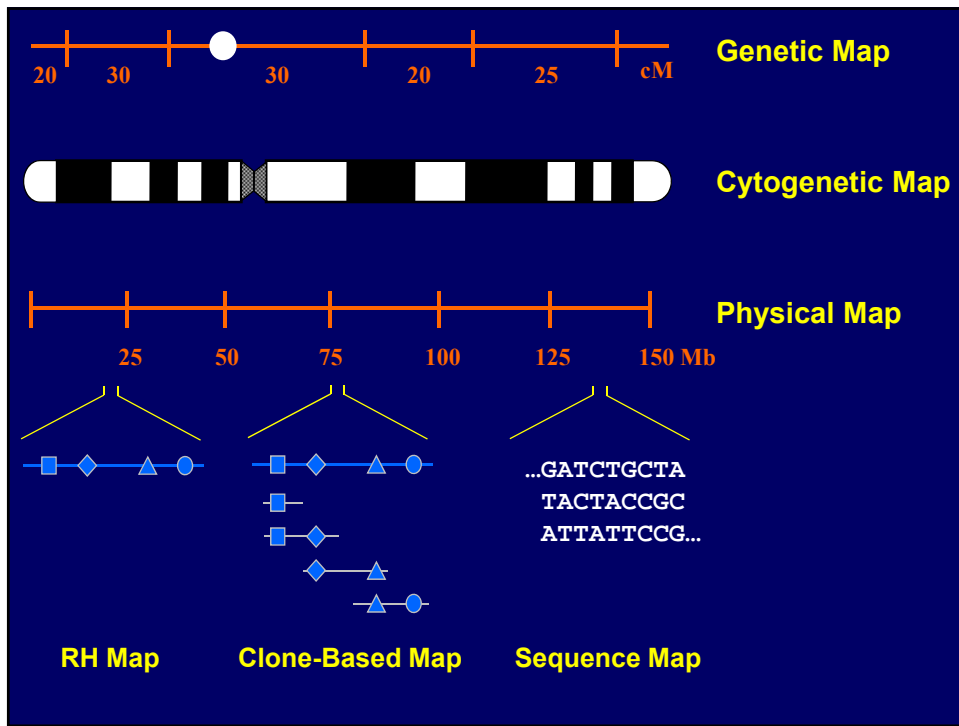


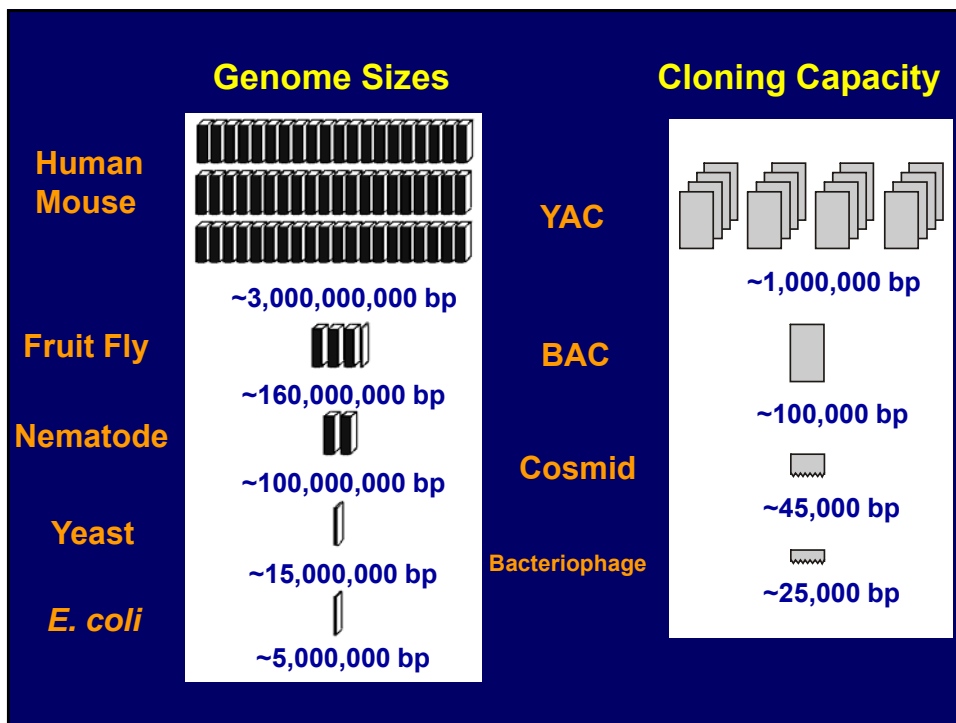
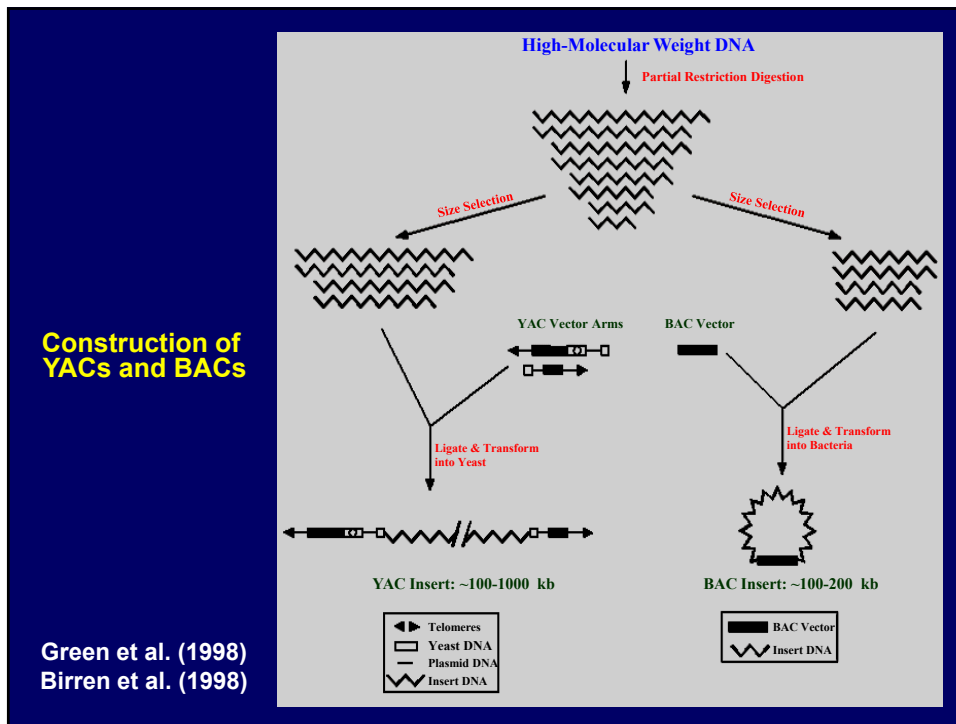


Outline

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





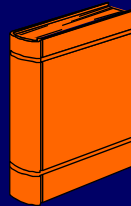


Bacterial Artificial Chromosomes (BACs)

- Bacterial-Based Cloning System
- Based on the *E. coli* F Factor (Fertility Plasmid): Replication Control
- Cloned Inserts: 100-200 kb, Circular DNA
- Low Copy Number
 - Low Yields of DNA by Standard Methods
 - Reasonably Stable
- See Birren et al. (1998)
- Availability of BAC Libraries from Many Vertebrate Species (e.g., www.chori.org/bacpac)



Genome
(~3000 Mb)



Chromosome
(~130 Mb)

G	A	T	C	T	C	T	A	G	A	A	T	C	T	C
G	A	G	A	T	C	T	C	T	A	G	A	T	C	T
G	T	G	G	A	C	T	T	T	G	A				
T	T	T	T	T	T	T	T	T	T	T	T	T	T	T
T	T	T	T	T	T	T	T	T	T	T	T	T	T	T
A	A	A	A	A	A	A	A	A	A	A	A	A	A	A
G	A	G	A	T	C	T	C	T	A	G	A	A	T	C
G	A	G	A	T	C	T	C	T	A	G	A	A	T	C
C	C	C	C	C	C	C	C	C	C	C	C	C	C	C
G	A	G	A	T	C	T	C	T	A	G	A	A	T	C
G	T	G	G	A	C	T	T	T	G	A				
G	T	G	G	A	C	T	T	T	G	A				

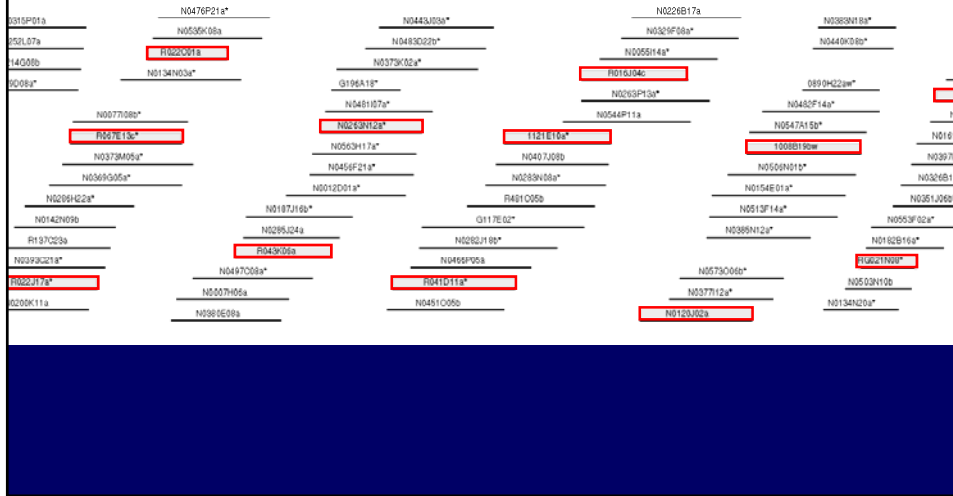
YAC
(~0.5-1.0 Mb)

G	A	T	C	T	C	T	A	G	A	A	T	C	T	C
G	A	G	A	T	C	T	C	T	A	G	A	A	T	C
G	T	G	G	A	C	T	T	T	G	A				
T	T	T	T	T	T	T	T	T	T	T	T	T	T	T
T	T	T	T	T	T	T	T	T	T	T	T	T	T	T
A	A	A	A	A	A	A	A	A	A	A	A	A	A	A
G	A	G	A	T	C	T	C	T	A	G	A	A	T	C
G	A	G	A	T	C	T	C	T	A	G	A	A	T	C
C	C	C	C	C	C	C	C	C	C	C	C	C	C	C
G	A	G	A	T	C	T	C	T	A	G	A	A	T	C
G	T	G	G	A	C	T	T	T	G	A				
G	T	G	G	A	C	T	T	T	G	A				

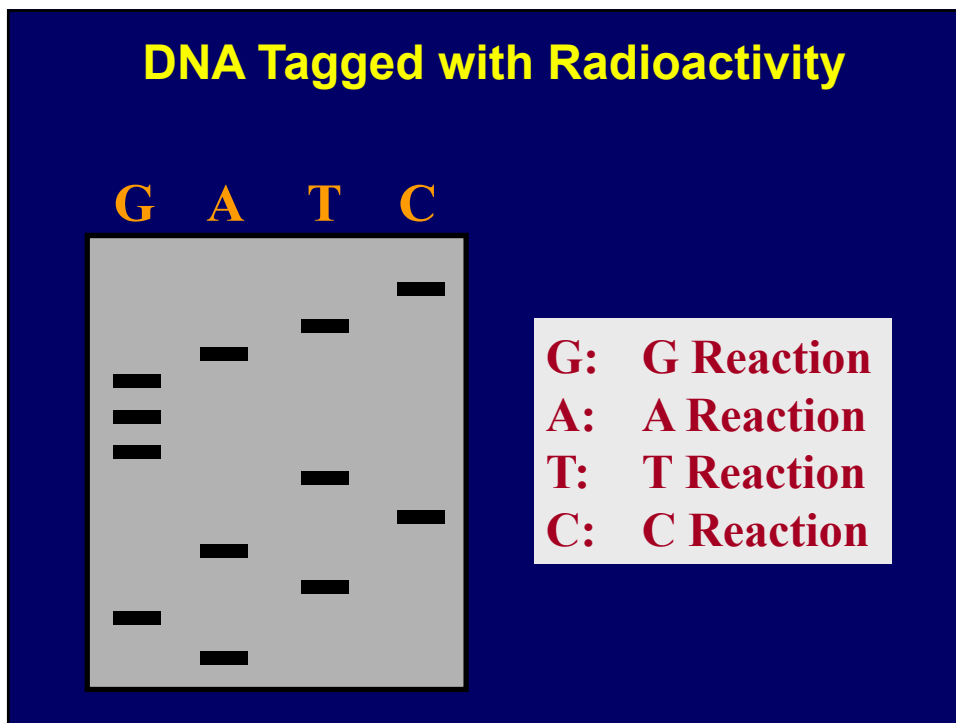
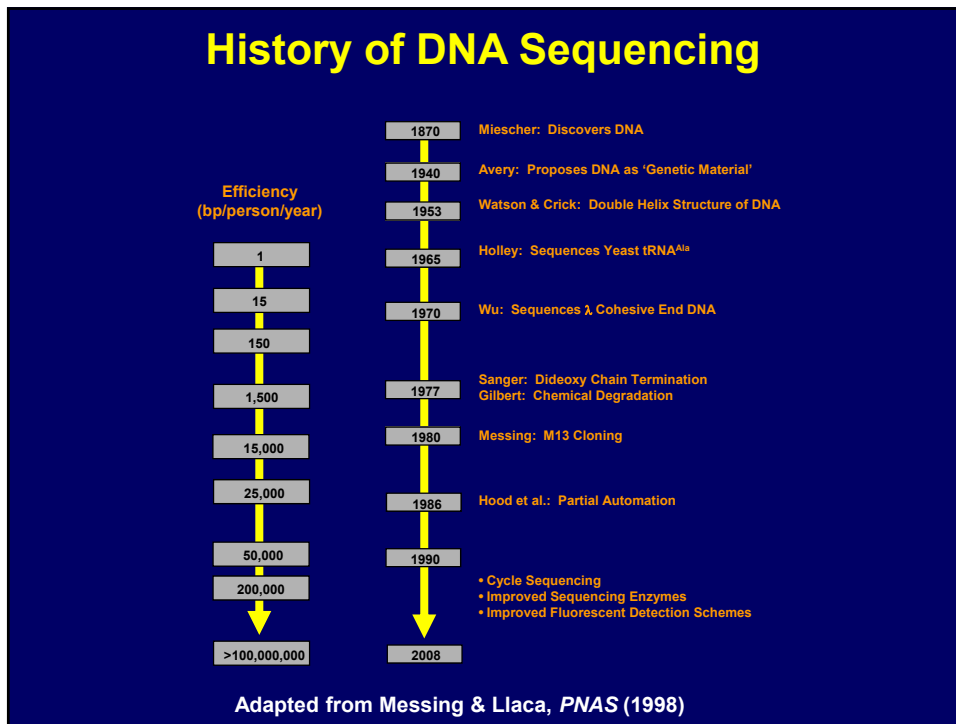
BAC
(~0.1-0.2 Mb)

Sequence-Ready BAC Contig Map

Marra et al. (1997)



DNA Sequencing

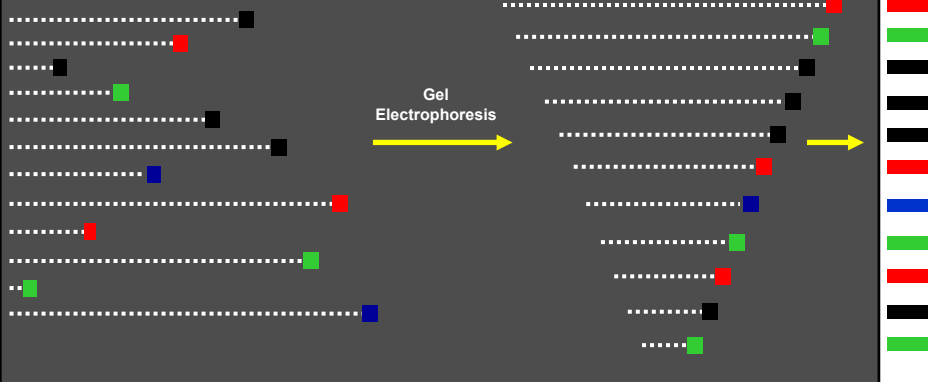


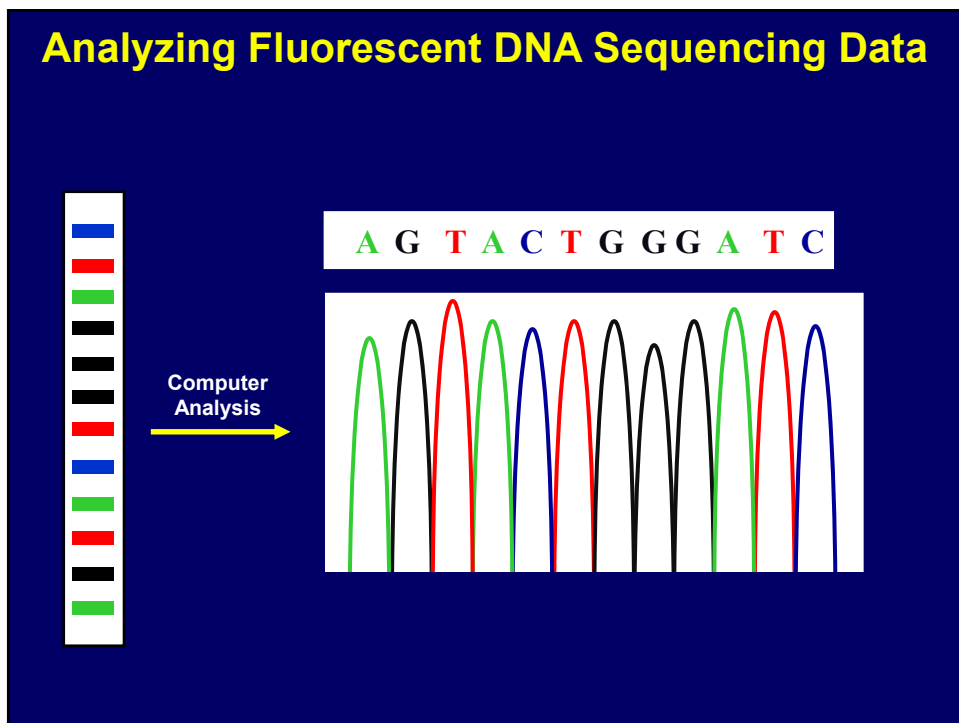
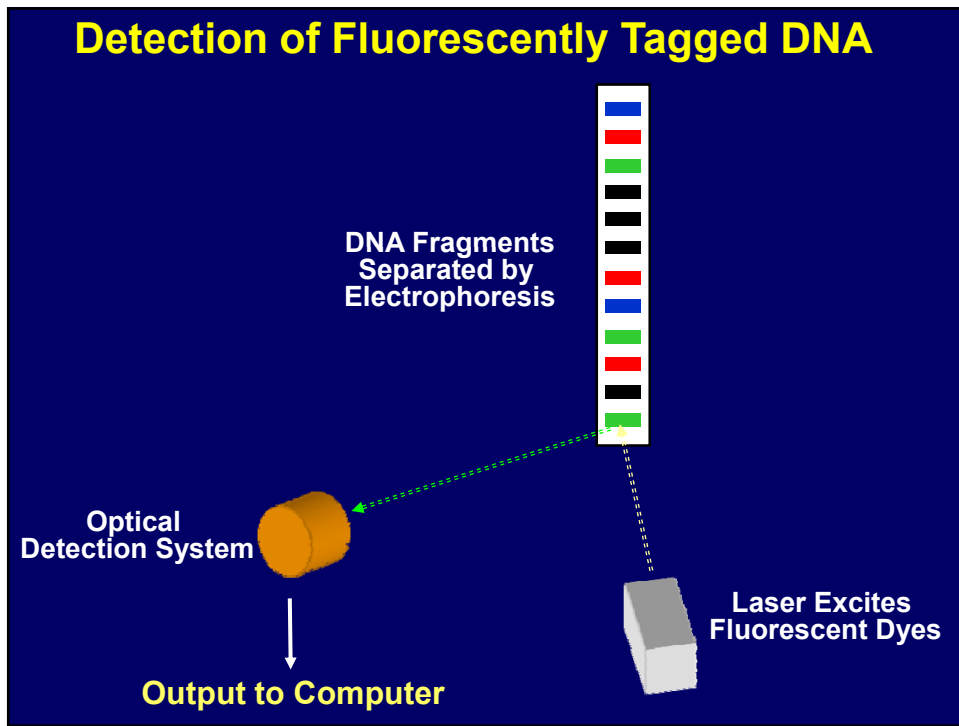
Radioactive Sequencing



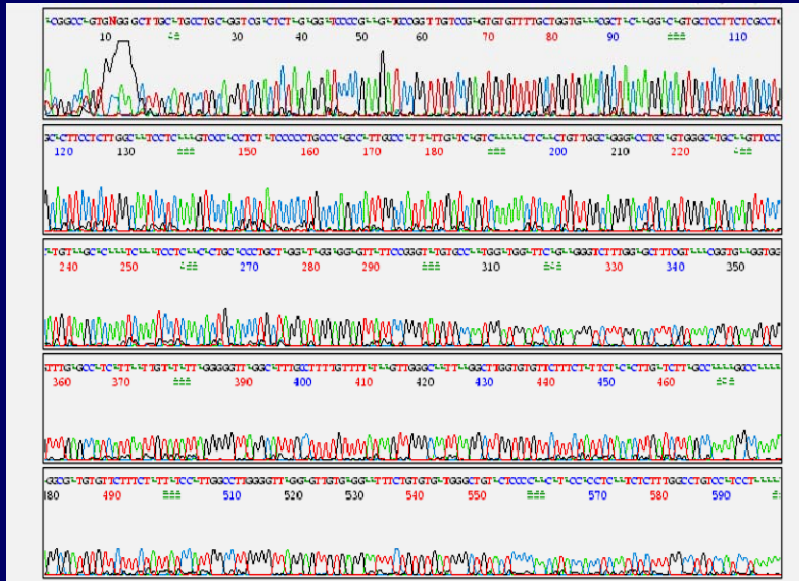
Fluorescent DNA Sequencing

AGTACTGGGATC





Fluorescent DNA Sequencing Results



Slab Gel-Based DNA Sequencing Instruments



Capillary-Based DNA Sequencing Instruments



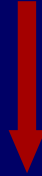
Large-Scale cDNA Sequencing

- 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)

Subclone Construction

```
GATCTGTAGATCTC
GAGTCTTGGAGTTC
GTGGGAACTGTGTA
TTTGACTGACAGAT
TACGTGTAGAGATG
ATGATGCACCTGACC
GGTTCGACTGTGAG
GACTCACTGCACCTCA
GAGGCGACGCGCGCT
GTGCACGTCGACACC
GATTTATACATTTA
AATCTTAGATTTAG
```

BAC DNA

Prepare Multiple Copies

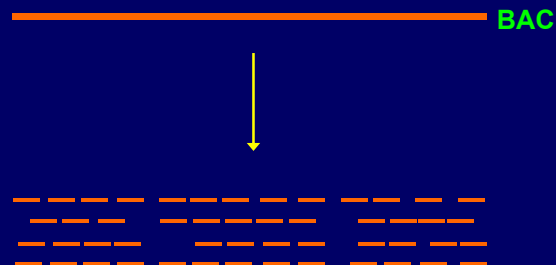
Randomly Fragment

Subclone Fragments

```
GA GA GA GATCTGTAGATCTC
GA GA GA GAGTCTTGGAGTTC
GA GA GA GTGGGAACTGTGTA
TT TT TT TTTGACTGACAGAT
AT AT AT ATGATGCACCTGACC
GA GA GA GGTTCGACTGTGAG
GA GA GA GACTCACTGCACCTCA
GA GA GA GAGGCGACGCGCGCT
GA GA GA GTGCACGTCGACACC
GA GA GA GATTTATACATTTA
AT AT AT AATCTTAGATTTAG
```



Shotgun Sequencing Strategy



Poisson Calculations

The sequencing strategy for the shotgun approach follows the Lander and Waterman application of the Poisson distribution

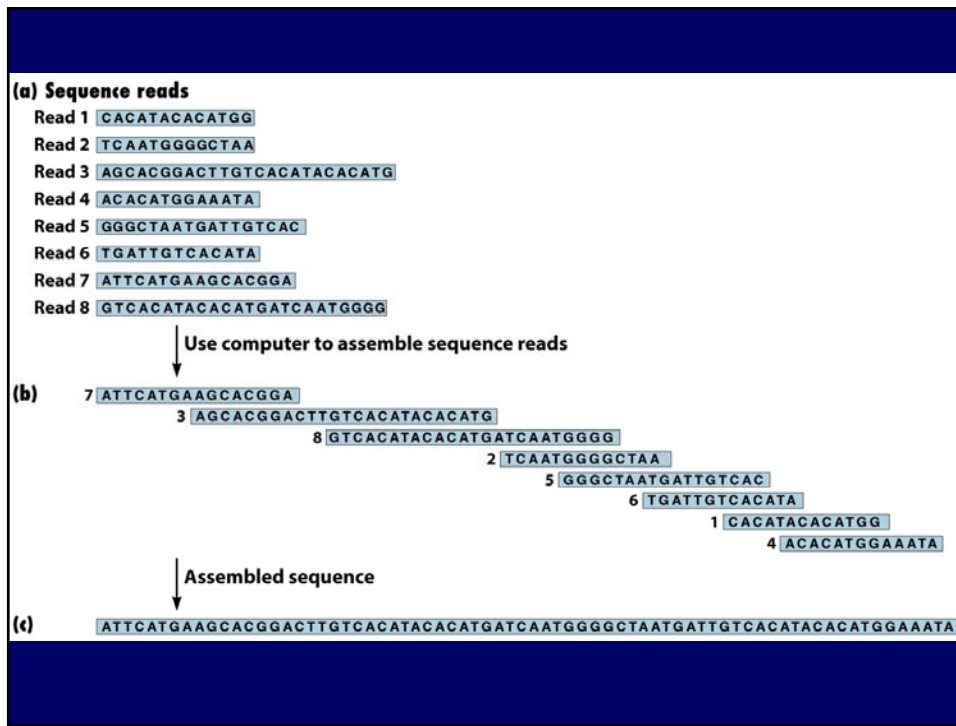
The probability a base is not sequenced is given by:

$$P_0 = e^{-c}$$

Where:

- c = fold sequence coverage ($c = LN/G$),
- LN = # bases sequenced, i.e. L = average sequencing read length and N = # reads
- G = target sequence length
- $e = 2.718$ ($e = 2.718281828459$)

Fold Coverage	$P_0 = e^{-c}$	% not sequenced	% sequenced
1	0.37	37%	63%
2	0.135	13.5%	87.5%
3	0.05	5%	95%
4	0.018	1.8%	98.2%
5	0.0067	0.6%	99.4%
6	0.0025	0.25%	99.75%
7	0.0009	0.09%	99.91%
8	0.0003	0.03%	99.97%
9	0.0001	0.01%	99.99%
10	0.000045	0.005%	99.995%



Shotgun Sequence Assembly

aligned reads

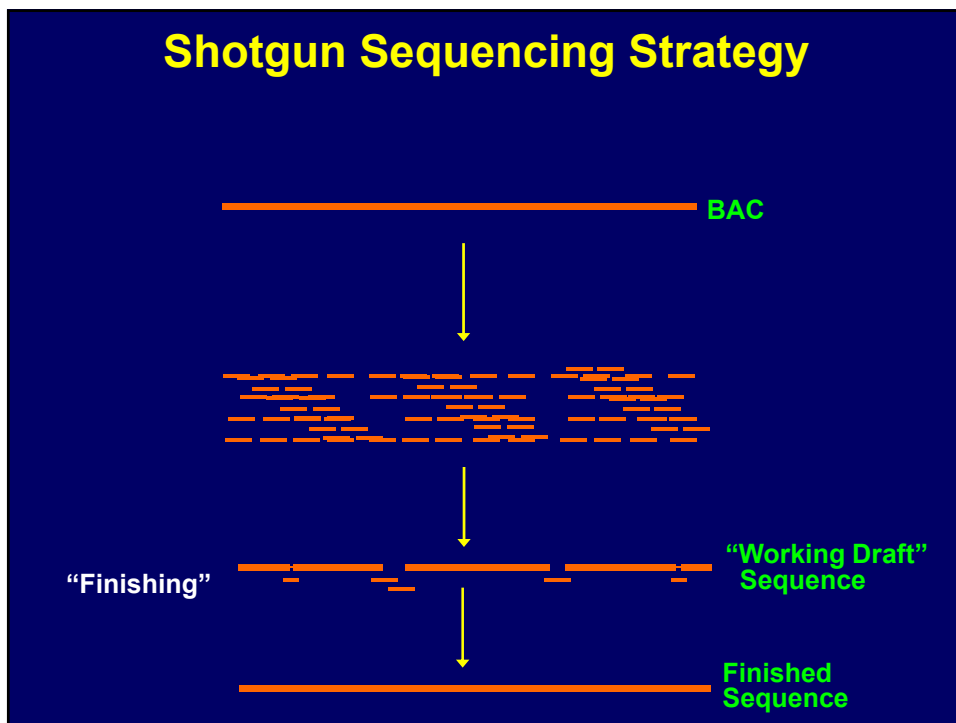
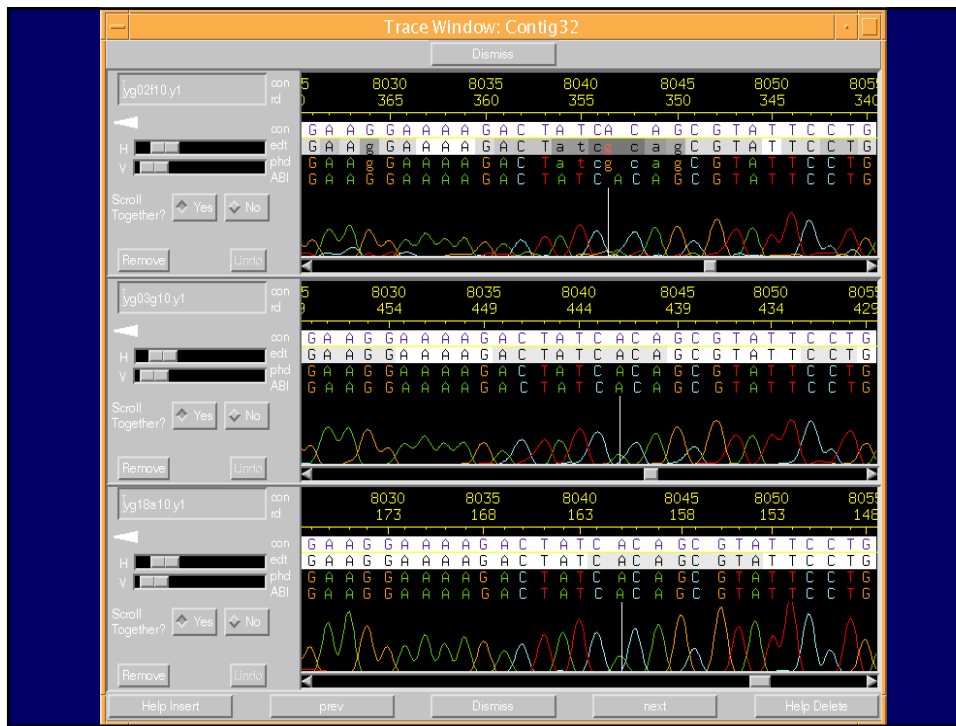
File Navigate Info Color Dim Misc Help

yg.fasta screen.a.ce.3 Contig32 Some Tags Pos

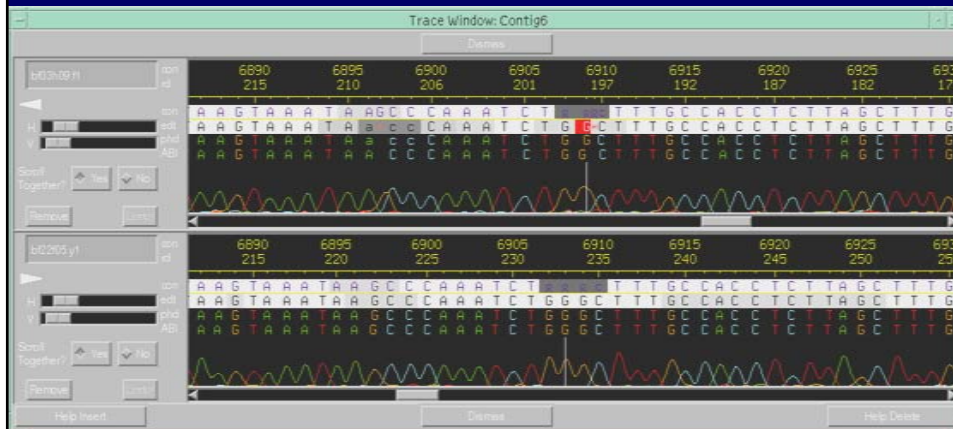
Search for String Compl Cont Compare Cont Find Main Win Exp Err/10kb 12.17

	8030	8040	8050	8060	8070	8080	8090	8100	811
CONSENSUS	AGGAAAAGACTATCACAGCGTATTCTGAAAGAGATGAACTATGAATTGAGTGTAGGCTTCTCTGCAGAGGCCAAR*GGTAGGAT								
yg12h02.x1	▶ cttgggggggaaagaaaacttttcccccgtttccctgaaggagaaacacactgaaatgggggggggttttttttgg								
yg03d09.y1	▶ aggaaaagactatcacagcgtaattcctgaaagagatgaaCTATGAattGAgTgtaggcttctctgcagaggcaaa*ggtaggat								
yg09g04.x1	▶ aggaaaagactatcacagcgtaattcctgaaagagatgaaactatgaaattgagtaggcttctctgcagaggcaaa*ggtaggat								
yg13h04.x1	▶ AGGAAAAGACTATCACAGCGTATTCTGAAAGAGATGAACTATGAATTGAGTGTAGGCTTCTCTGCAGAGGCCAAR*GGTAGGAT								
yg01e03.y1	▶ AGGAAAAGACTATCACAGCGTATTCTGAAAGAGATGAACTATGAATTGAGTGTAGGCTTCTCTGCAGAGGCCAAR*GGTAGGAT								
yg08h10.x1	▶ xxxxxxxxxxxxatcacagcgtaattcctgaaagagatgaaactatgaaattgagtaggcttctctgcagaggcaaa*ggtaggat								
yg04f11.y1	▶ xxxxxxxxxxxxxxxxxxxxxxxxxxxcagctcgccca								
yg01g01.y1	▶ acatcgttcaagttgaacatccgctatxx								
yg01g07.y1	▶ xxx								
yg02e04.y1	▶ AGGAAA								
ygGAAAAGACTATCACAGCGTATTCTGAAAGAGATGAACTATGAATTGAGTGTAGGCTTCTCTGCAGAGGCCAAR*GGTAGGAT									
yg02f10.y1	▶ aggaaaagactatcacagcgtaattcctgaaagagatgaaactatgaaattgagtaggcttctctgcagaggcaaa*ggtaggat								
yg02c10.y1	▶ AGGAAAAGACTATCACAGCGTATTCTGAAAGAGATGAACTATGAATTGAGTGTAGGCTTCTCTGCAGAGGCCAAR*GGTAGGAT								
yg18a10.y1	▶ AGGAAAAGACTATCACAGCGTATTCTGAAAGAGATGAACTATGAATTGAGTGTAGGCTTCTCTGCAGAGGCCAAR*GGTAGGAT								
yg03g10.y1	▶ aggaaaagactatcacagcgtaattcctgaaagagatgaaactatgaaattgagtaggcttctctgcagaggcaaa*ggtaggat								
yg08f02.y1	▶ aggaaaagactatcacagcgtaattcctgaaagagatgaaactatgaaattgagtaggcttctctgcagaggcaaa*ggtaggat								
yg02h10.y1	▶ agaaaaatcctatccagcgtaattcctgaaagagatgaaactatgaaattgagtaggcttctctgcagaggcaaa*ggtaggat								
yg18e09.y1	▶ aggaaaagactatcacagcgtaattcctgaaagagatgaaactatgaaattgagtaggcttctctgcagaggcaaa*ggtaggat								
yg13d05.y1	▶ ggtgcgcggtcaactgtgcccggtctgcgtgcgccga*tgccgcgc								

“Consed” (Gordon et al., 1998)

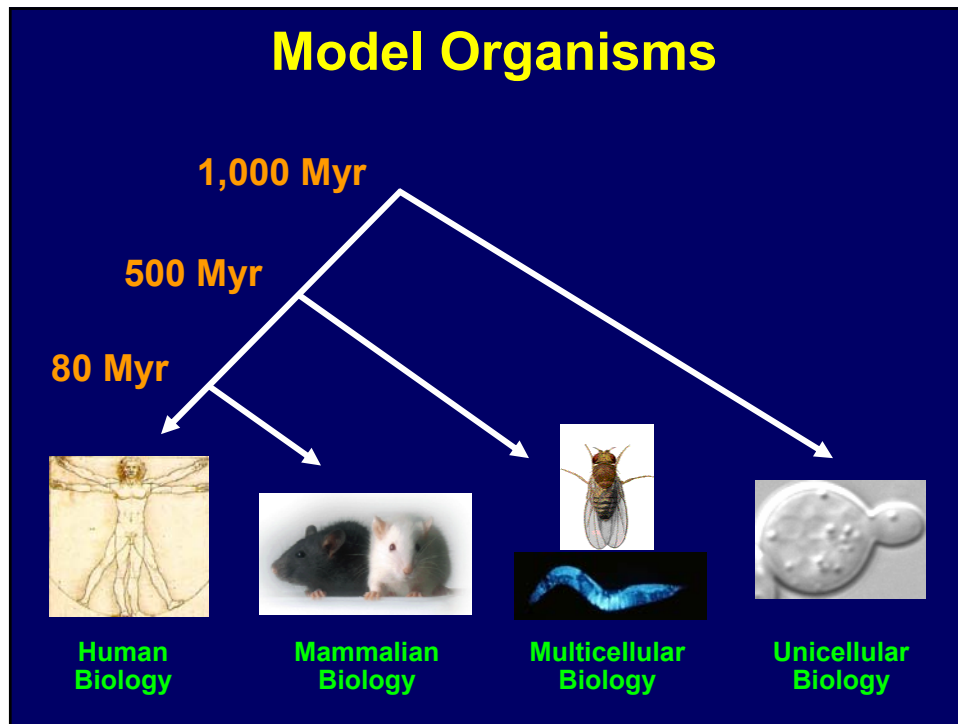


Sequence Finishing: Resolving Ambiguities




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

Historically Significant Genome Sequencing Projects



Microbial Genome Sequences



Comprehensive Microbial Resource

Search Locus for

Genome Search

Organism name:

Genome List

Gene Search

Search by:

Locus:

Match: Exact Inexact

Keywords/Accession:

Data Summary

	Complete	Draft	Totals
Bacteria	353	17	370
Archaea	28	0	28
Viruses	3	0	3
Totals	384	17	401

Welcome to the Comprehensive Microbial Resource

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. A [CMR Mirror](#) site maintained by the Genome Encyclopedia of Microbes (GEM) in Korea is also available. [\[More Information\]](#) [\[Publication Information\]](#)

CMR Menu Bar Tools

CMR offers a wide variety of tools and resources, all of which are available off of our menu bar at the top of each page. Below is an explanation and link for each of these menu options. First time users can use our [CMR tutorial](#) to learn how to navigate this site.


Genome Tools

Find organism lists as well as summary information and analyses for selected genomes.

Searches

Search CMR for genes, genomes, sequence regions, and evidence.

Announcements



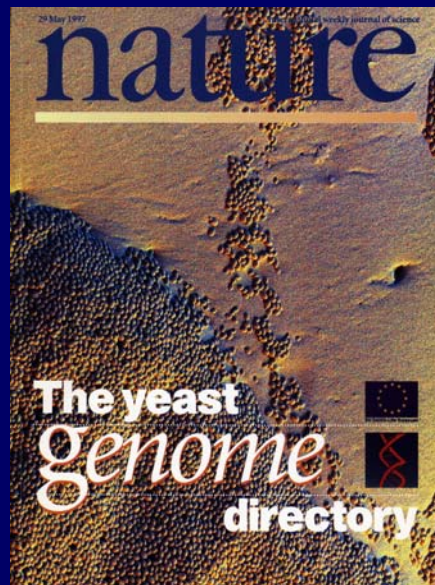
March 13, 2007: CAMERA is a web resource for metagenomic research. CAMERA's debut coincides with the [publication](#) of the [Global Ocean Sampling](#) expedition's extensive dataset cataloging over 6 million new genes from uncultured marine microbes. Come visit [CAMERA](#), and see our growing collection of metagenomics datasets and tools.

Latest Releases

Data Release: 21.0

www.tigr.org

First Eukaryotic Genome Sequence



Goffeau et al. (1997)

First Animal Genome Sequence

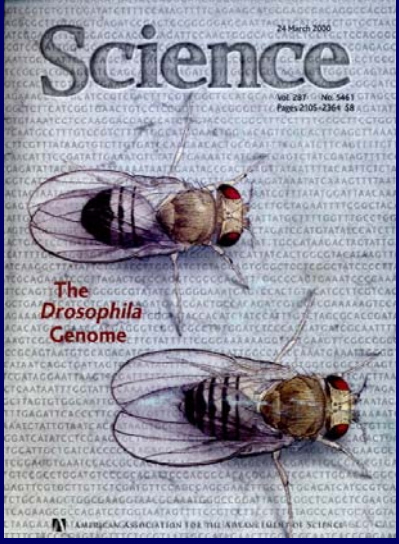


Genome Sequence of the Nematode *C. elegans*:
A Platform for Investigating Biology

The *C. elegans* Sequencing Consortium*

C. elegans Sequencing Consortium (1998)

Second Animal Genome Sequence



The Genome Sequence of *Drosophila melanogaster*

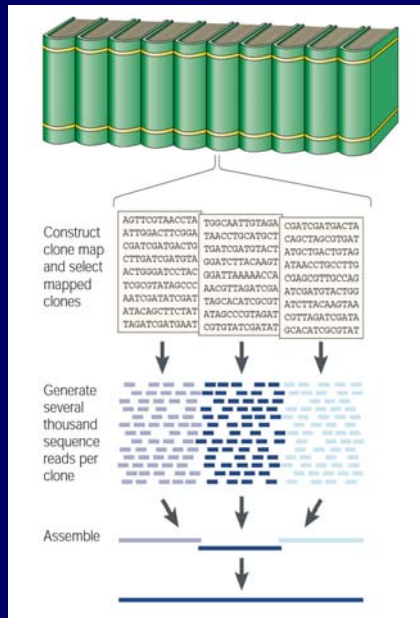
REVIEW

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,² Peter W. Li,¹ Roger A. Hoskins,² Richard F. Galle,¹ Reed A. George,² Suzanna E. Lewis,¹ Stephen Richards,² Michael Ashburner,¹ Scott N. Henderson,¹ Granger G. Sutton,¹ Jennifer R. Wortman,¹ Mark D. Vandell,¹ Qing Zhang,¹ Jin Xi Chen,¹ Rhonda C. Brandon,¹ Yu-Hui C. Rogers,¹ Robert G. Blasziak,² Mark Champe,² Berret D. Pfeiffer,² Kenneth H. Wan,² Clare Doyle,² Evan G. Baxter,² Gregg Helt,² Catherine R. Nelson,² George L. Gabor Miklos,² Josep F. Abril,² Anna Aghayani,² Hui-Jin An,² Cynthia Andrews-Pfannkoch,¹ Daria Baldwin,¹ Richard M. Bailey,¹ Anand Baner,¹ James Baxterdale,¹ Leyla Bayraktaroglu,² Ellen M. Beasley,² Karen Y. Beeson,² P. V. Benos,¹⁸ Benjamin P. Berman,² Deepali Bhandari,¹ Slava Bolshakov,¹¹ Dana Borkova,² Michael R. Botchan,¹³ John Bouck,² Peter Brokstein,² Philippe Brotier,¹⁴ Kenneth C. Burtis,¹⁵ Dana A. Busam,¹ Heather Butler,¹⁶ Edouard 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,² Anne Deslattes Mays,¹ Ian Dew,¹ Suzanne M. Dietz,¹ Kristina Dodson,¹ Lisa E. Doup,¹ Michael Downes,²¹ Shannon Dugan-Rocha,² Boris C. Dunkov,²² Patrick Dunn,²³ Kenneth J. Durbin,²⁴ Carlos C. Evangelista,¹ Concepcion Ferraz,²⁵ Steven Ferrieri,¹ Wolfgang Fleischmann,¹ Carl Foster,¹ Andrei E. Gabrielian,¹ Neha S. Garg,¹ William M. Gelbart,² Ken Glasser,¹ Anna Glöckl,¹ Fangcheng Gong,¹ J. Harley Gorrell,² Zhiping Gu,¹ Ping Guan,¹ Michael Harris,¹ Nomi L. Harris,² Damon Harvey,¹ Thomas J. Heiman,¹ Judith R. Hernandez,² Jarrett Houck,² Damon Houston,¹ Kathryn A. Houston,² Timothy J. Howland,¹ Ming-Hui Wei,¹ Chinyere Ibegwam,¹ Mena Jalali,¹ Francis Kalish,¹ Gary H. Karpen,¹¹ Zhaod Ke,¹ James A. Kennison,²⁴ Karen A. Ketchum,¹ Bruce E. Kimmel,² Chinnappa D. Kodira,¹ Cheryl Kraft,¹ Saul Kravitz,² David Kulp,²⁶ Zhongyu Lei,¹ Paul Lasko,²⁷ Yiding Lei,¹ Alexander A. Levitsky,¹ Jayin Li,¹ Zhenyu Li,¹ Yong Liang,¹ Xiaoying Lin,²⁸ Xiangjun Liu,¹ Bettina H. Matesa,¹ Tina C. McIntosh,¹ Michael P. McLeod,² Duncan McPherson,¹ Genady Mefkovic,¹ Natalia V. Milshina,¹ Clark Moberly,¹ Joe Morris,¹ Ali Moshrefi,² Stephen M. Mount,²⁷ Mae Moy,¹ Brian Murphy,¹ Lee Murphy,²⁹ Donna M. Muzny,² David L. Nelson,² David R. Nelson,²⁹ Keith A. Nelson,² Katherine Nixon,¹ Deborah R. Nusbaum,¹ Joanne M. Paclik,¹ Michael Palazzolo,³ Gijung S. Pitman,¹ Sun Pan,¹ John Pollard,¹ Vinita Puri,¹ Martin G. Reese,¹ Knut Reineert,¹ Karin Remington,¹ Robert D. C. Saunders,²⁰ Frederick Schaefer,¹ Hua Shen,¹ Bizhang Christopher Shue,¹ Inga Sidén-Kiamos,¹ Michael Simpson,¹ Marian P. Skupski,¹ Tom Smith,¹ Eugene Spier,¹ Allan C. Spradling,¹ Mark Stapleton,² Renee Strong,¹ Eric Sun,¹ Robert Svikas,²⁶ Cyndee Tector,¹ Russell Turner,¹ Eli Venter,¹ Alhui H. Wang,¹ Xin Wang,¹ Zhen-Yuan Wang,¹ David A. Wasserman,²⁸ George M. Weinstock,² Jean Weissenbach,¹ Sherita M. Williams,¹ Trevor Woodage,¹ Kim C. Worley,² David Wu,¹ Song Yang,¹ Q. Allison Yao,¹ Jane Ye,¹ Hu-fang Yeh,¹ Jayshree S. Zaveri,¹ Hing Zhan,¹ Guangren Zhang,¹ Qi Zhao,¹ Liangsheng Zheng,¹ Xiangjun H. Zheng,¹ Fei Ni Zhong,¹ Wuyuan Zhong,¹ Xiaojin Zhou,¹ Shileiping Zhu,¹ Xiaohong Zhu,¹ Hamilton O. Smith,¹ Richard A. Gibbs,¹ Eugene W. Myers,¹ Gerald M. Rubin,²⁴ J. Craig Venter¹

Adams et al. (2000)

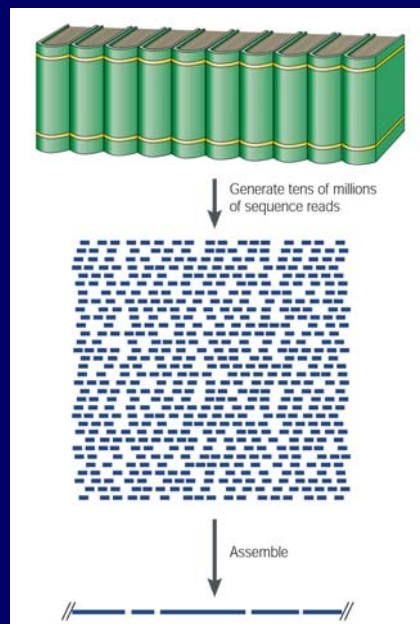


Clone-Based Shotgun Sequencing



Green (2001)

Whole-Genome Shotgun Sequencing

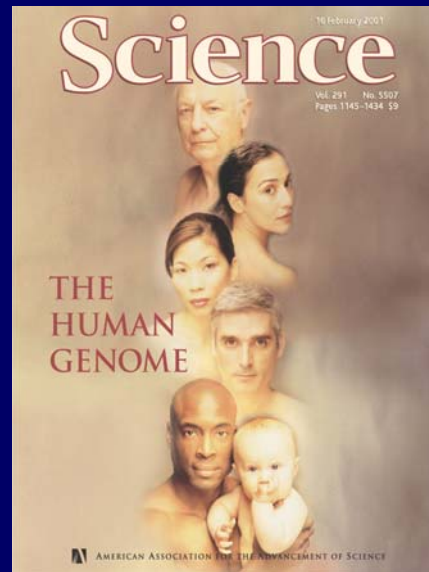


Green (2001)

February, 2001 Draft Sequence

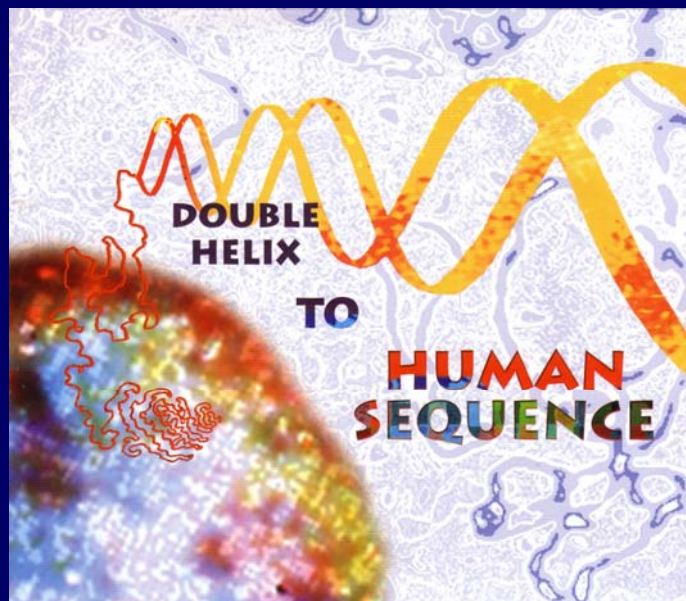


International Human Genome Sequencing Consortium (2001)



Venter et al. (2001)

April, 2003 Completion



October, 2004 Publication



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.

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 status of this finishing process. The draft genome sequence (total 283 million base-pair contigs) interrupted by only 2,811 gaps, it covers ~99% of the nucleotide sequence and is accurate to an error rate of ~1 error per 100,000 bases. Many of the remaining nucleotide gaps are associated with repetitive sequences 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 history. The human genome sequence is available 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 high-resolution sequence of the non-repetitive portion of the human genome. The initial work followed a new conceptual approach: (1) the mapping of the human genome, (2) the sequencing of the human genome, and (3) the assembly of the genome into a high-resolution map. The sequencing of the genome was a challenge because of the large size of the genome, the high degree of sequence similarity between human chromosomes, and the presence of repetitive DNA. The International Human Genome Sequencing Consortium (IHGSC), an open collaboration involving twenty centers in six countries, was created to coordinate the sequencing of the human genome. In February 2001, the IHGSC and Celera Genomics* each reported a preliminary draft sequence of the human genome. These sequences allowed systematic study of the human genome, including identification of genes, cellular and molecular architecture of genes, regional differences in genome complexity, distribution and history of transposable elements, distribution of polymorphisms and relationship between genetic, biochemical and physical distance. Moreover, systematic knowledge of the human genome has enabled new tools and approaches that have revolutionized biomedical research.

Both draft sequences, however, had important shortcomings. The IHGSC sequence, for example, covered ~90% of the euchromatic genome. It was interrupted by ~130,000 gaps and the order and orientation of many segments within local regions had not been established. The Celera sequence, in contrast, covered the euchromatic genome, but was interrupted by ~1.5 million gaps. The Celera sequence was defined as having an error rate of, at most, one error per 100 bases, and the gaps in its sequence were coverage in finished segments of at least 90% of the euchromatic genome, with the only gaps being those relevant to all available technologies (see <http://www.genome.gov/10899723>). The goal was challenging because the human genome is complex with nucleosomes, repetitive DNA and large segmental duplications, which greatly complicate the determination of genome structure. In fact, some complete sequences have been revealed to be, in some respects, more complex than the genome. The genome is, in fact, more complex than the genome and has much smaller structure.

We conclude from the results of a mid-year effort by the IHGSC towards the goal of a complete human genome. The number of gaps has been reduced 40-fold to only 281, most of which are associated with segmental duplications and will require new methods for resolution. The assembled near-complete genome has an error rate of only ~1 error per 100,000 bases. It contains 2,811 gaps, nucleotide and sequence ~99% of the nucleotide genome. This paper describes the current genome sequence and the process used to produce it, discusses the accuracy and completeness of the sequence and illustrates biological insights made possible by the sequence. We also discuss how a complete human genome of the order of the human genome, the initial draft sequence was previously reported* and a series of papers is being written describing the technical achievements, the initial organization of genes and other features.

Current genome sequence

Human genome

The process of converting the initial draft sequence into a near-complete sequence is defined as 'finishing'. It is a complex, iterative process that generally involves multiple rounds of sequencing, mapping, and assembly. The finishing process is a complex, iterative process that generally involves multiple rounds of sequencing, mapping, and assembly. The finishing process is a complex, iterative process that generally involves multiple rounds of sequencing, mapping, and assembly. The finishing process is a complex, iterative process that generally involves multiple rounds of sequencing, mapping, and assembly.

Finally, the finishing process involved two distinct components: (1) producing finished maps, consisting of continuous and accurate paths of overlapping large insert clones spanning the euchromatic region of each chromosome arm and (2) producing finished clones, consisting of continuous and accurate nucleotide sequence across each large insert clone. In practice, these two components were tightly interrelated in that progress in each often depended on results from the other. The sequence was finished by August 2001. Further information about the finishing process and finishing methods can be found in the Supplementary Information (S1) and at <http://www.genome.gov/10899723>.

In total, we generated a genome sequence from 30,208 large-insert clones (total length ~1.84 gigabases (Gb)) and finished the genome from 45,742 of these clones (total length ~1.67 Gb). The clones consisted primarily of bacterial artificial chromosomes.

984

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

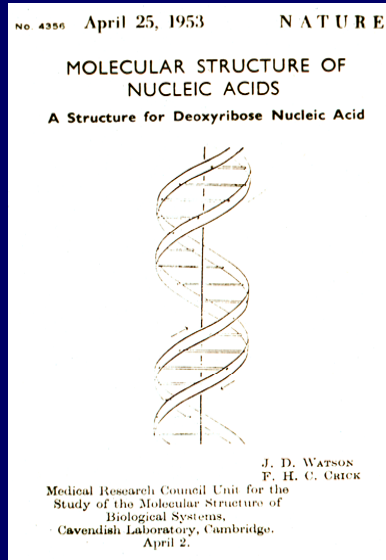
Human genome mapping ranks No. 1 in health news

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

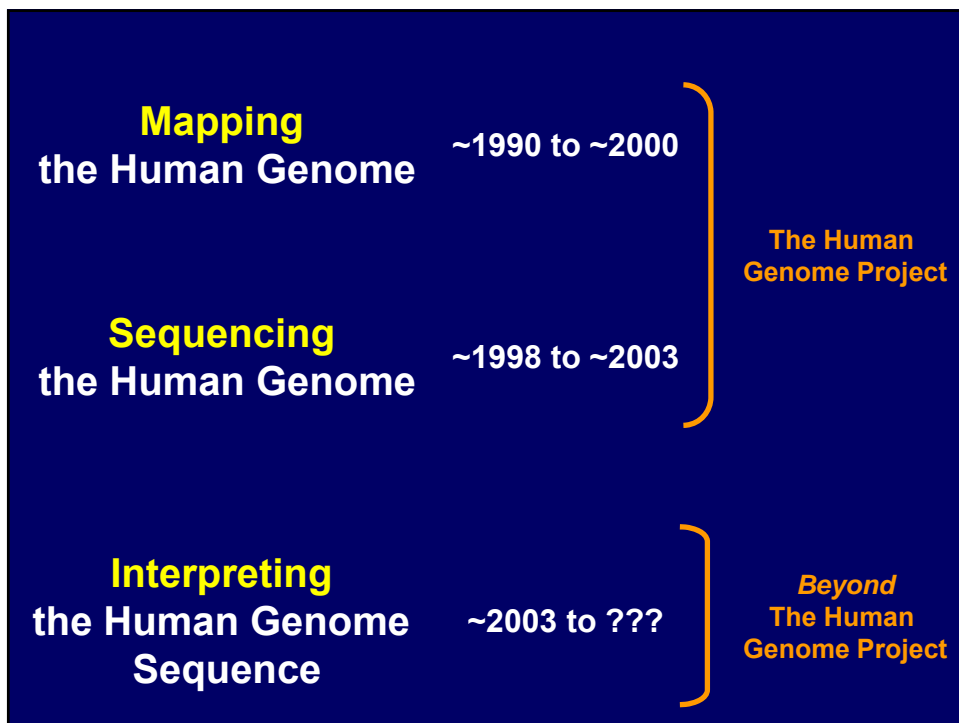
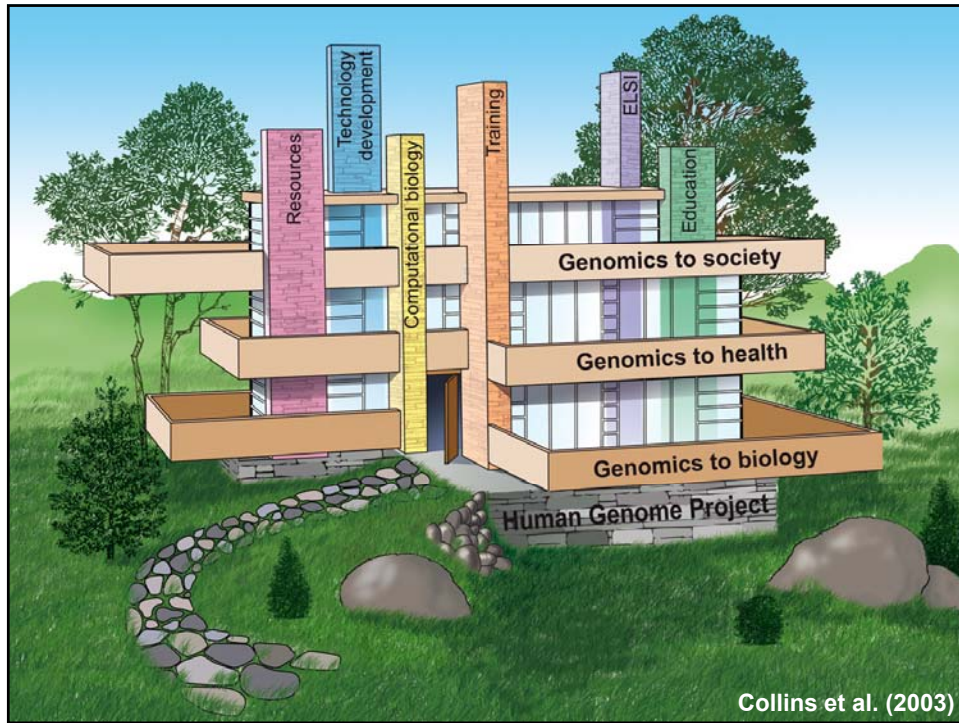
April, 1953 → **April, 2003**



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

What's Next?







~3,000 bp (0.0001%) of Human Genome Sequence

TCGCCGGAACTTTCCGGCTCTCTAAGGCTGTATTTTGTATATACGAAAAGGCACATTTCCCTCCCTTTTCAAATGCACCTTGCAAACGTAACAG
 GAACCCGACTAGGATCAFCGGGAAAAGGAGGAGGAGGAAGGCAGGCTCCGGGGAAGCTGGTGGCAGCGGGTCTGGGTCTGGCGGACCCCTGA
 CGCGAAGGAGGGCTAGGAAGCTCTCCGGGGAGCCGGTTCTCCCGCGGTGGCTTCTCTGTCTCCAGCCTGGCAACTGGACCTAAAGAGAGG
 CCGCGACTGTCCGCCACCTGCGGGATGGGCTGGTCTGGGCGGTAAGGACACGGACCTGGAAGGAGCGCGCGAGGGAGGGAGGCTGGGAGTC
 AGAATCGGAAAAGGAGGTGCGGGCGCGGAGGAGCGAAGGAGGAGGAGGAAGGAGCGGGAGGGGTGCTGGCGGGGTGCGTAGTGGGTGGA
 GAAAGCCCTAGAGCAAATTTGGGGCCGACCCAGGCAGCACTCGGCTTTAACTGGGCACTGAAGGGGGGAAAGAGCAAAGGAAAGGGGTGG
 TGTGCGGAGTAGGGTGGTGGGGGAAATGGGAAGCAAATGACATCACAGCAGGTCAGAGAAAAGGGTTGAGCGGCAGGCCACCCAGGATAGTAG
 STCTTTGGCATTAGGAGCTTGAAGCCAGACCGCCCTAGCAGGGACCCAGCGCCGAGAGACCAATGCAGAGGTCGCCCTCTGAAAAGGCCAGCT
 TGTCTCAAACCTTTTTTTCAGGTGAGAAGGTGGCCACCAGGCTTCGAAAAGCACGTCGCCACGAAAGAGGGCGTGTGTATGGTGGGTT
 TGGGTAAGGAATAAGCAGTTTTTAAAAAGATGCGCTATCATTGTTTTGAAAGAAAATGGGATATTAGATAAAAACAGAAAGCATT
 AGAAGAGATGGAAGAATGAAGTGAAGCTGATTGAATAGAGAGCCACATCTACTGCAACTGAAAAGTTAGAACTCAAGACTCAAGTACGCTACT
 ATGCACTGTTTTTATTTTCAAGAACTAAAAAATCTTGTAAATAGTACCTAAGTATGGTTATGGTTTTCCCCCTCATGCCTGG
 ACACCTGATGCTCTTGGCACATACAGGTGCCATAGCTGCATATAGTAAAGTGCCTCAGAAAACATTTCTTGACTGAATTCAGCCAAACAAAAT
 TTGGGTTAGGTAGAAAATATGCTTAAAGTATTTTGTATGAGACTGGATATATCTAGTATTTGTCACAGGTAATGATCTTCAAATAATG
 AAAGCAAATTTGTGAATATTTTAAAAAGTTACTTCACAAGCTATAAATTTAAAAAGCCATAGGAATAGATACCGAAGTTATATCCAA
 CTGACATTTAATAAATTTGATTTTCATAGCCTAATGTGATGAGCCACAGAAGCTTGCAAACCTTAATGAGATTTTTAAAAATAGCATCAAGTTCGG
 AACTTAGGCAAAGTGTGTAGATGAGCACTCATATTTGAAGTGTCTTGGATATGCACTACTTTGTTCTGTTATATAGTGGTGTGA
 ATGAATGAATAGTACTGCTCTCTTGGACATTACTTGAACATAATTAACCAATGAATAGCATACTGAGGTATCAAAAAGTCAAATATGT
 TATAAATAGCTCATATATGTGTAGGGGGAAAGGAATTTAGCTTTTCACTCTCTTATGTTTTAGTTCTCTGCATGTGCAGTTAATCCTGGAAC
 TCCGGTCAAGGAGAGACTGTTGGCCCTTGAAGGAGAGCTCCTCCCTGTGGATGAGAGAGAAGGACTTTACTCTTTGGAATATCTTTTTGTGT
 TGATTTATCCACTTTTGTACTCCACTATAAATCGGCTTACTATTGATCTGTTTTCCCTAGTCTTATAAAGTCAAAAATGTAATGGCAT
 AAATATAGACTTTTTTAGCAGAGAATTTGAGGAACCTAAATGCAACAGTCTFAAAAATGCACTTTTCAGAAAGATGAATATTTCACTGGATA
 GTTCTAAATCTAATGAACCTTAAAAATAGCTTACTATGATCTGTCAAAAGTGGGTTTTTATAAATTTCTTTTACAAAATCACCGACACATTT
 AATATAGTTAAAAATGCTATCAGGCTGGTTTGCAAAAGAAAATGTATTCAAAGGCTGCTAAGTGTGTTAAGGACATCTCATTCTGTTCTCC
 AAAATATTTCATAAGGTGCTTTAAGAAATAGGTATGTTTTTAAAAAGTTAAGTTCCTACTATTTATAGGAACTGACAAATCACCTAAAAATA
 CCAATGA
 TTACAAACTTCCTTCTGGCCTCTCGACTGCAATCTAAAAGTGTAAAAAACATATTTCTGCATTAAAGTGGCAGTATTGTTTACAA
 GTGGTAGGCTTGGAGTCAGATTTTGTATTGATCAGATCCTACATCTACTGTTAGTAGCTGTTGCTGAGGCAGGTCCTTAACATCTCTGTG
 TGTCACTGACCTTAAAAATTTGGAGACTCTCATAGGGGTTAATGCGCTTGAGAAAATGAAATGTGAAAAGTTAGCCCTAATGTAACTGCTATT
 ATGGATTACCATATTTTCACTTCACTCACAGTACATGCACCTTGTAAATAAAGATGCTCAATTCATCTTTGAGTATAATTTTGTGACTCTCAAT
 CTGGATATGCAATGAGTGGCCTGTATGAGAATTAATTTATGAAAATTTGTTTTCACATGGCCTTACCAGATATACAGGAACACGCTCACATG
 TTTCTATGTTGTTAAATGCCTTAGAATTTAACTTCTGAATAGGATCCCTCAGTTTGAGAGTCATAAAAAGAGTAAATATTATGGTAT

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

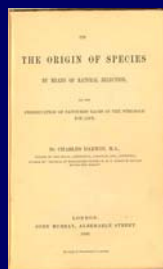
Gene Regulatory Elements

Chromosomal Functional Elements

Undiscovered Functional Elements (NOT Yet in Textbooks!)

Poor Inventory at Present

Foundational Milestones in Genetics & Genomics



Darwin

1859



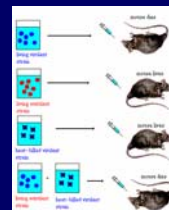
Mendel

1865



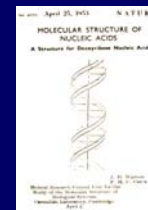
Miescher

1871



Avery

1944



Watson & Crick

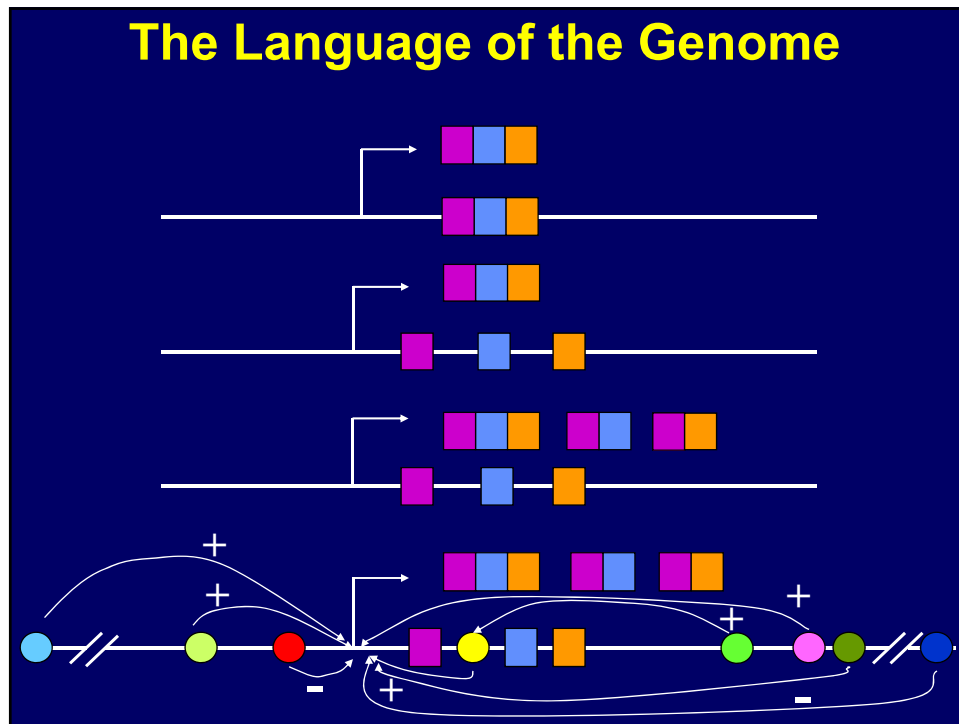
1953

Comparing Genomes is Like Cryptography

CKQEBHEREYTWASUISZMEISDFOGETHEBLPBGODFQSTLKSTUFFRTAC
DLUCEHEREZBRTTOISAWNCDARJJPTHERROFGODERGHCLSTUFFBRHA

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



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

Species A

```
GATCGTCTAGAATCTCGAGATC
TCTGAGAGTCTGGGAAACTGT
GTGATGTGACGATTTAGCCACA
GTTACGTFGAGAGATGATGA
TGCACCTGACCCGGTTTCACT
CTCAACGACTCACTCCACCTCA
GAGGCCACCCGCGCTGTGCAC
TACCGAGATACAGATACTTAC
ACAGGTTGTGACACCCCTTACC
CGTCCACACAGACTCACTCC
ACCTCAGAGGCCACCCGCGCT
GTGCACCTACCGAGATACAGAT
ACCTACACAGGTTGTGACACAG
ATCCTTACACAGCTTACAGATT
ACCATATATCCACTTACACAC
ATACCTACCCATTGACACCT
ATTATTATTACCGGACCGAGG
```

Compare


Species B

```
TATCGGCTAGAATCTCGAGATC
TCTGAGAGTCTGGGAAACTGT
GTGATGTGACTAGCCACAGTTA
CGTGTGAGAGATGATGATGCA
CCTGACCCGGTTTCACTCTCA
ACGACTCACTCCACTCAGAGG
CCACCGCGCTGTGCACGTCC
ACCAAGATOUTTACCACACTTA
CACATCACTCTCAGAGCTCAC
TCCACTCAGAGGCCACCCGCG
GCTGTGCACTCCACACGATC
CTTACCACCTTACACATTACC
ATATATCCACTTACCACACATA
CCTTACCAGATATCCACTACC
ACCATATATCCACTTACACAC
ACCTATTATTATTACCGAGGGA
GAGGGGTGACCACACTGTGACA
```


Sequences in Common (i.e., 'Conserved' or 'Constrained')

```
GATCGTCTAGAATCTCGAGATC
CTTGGAGTCTGGGAAACTGT
CGTGTGAGAGATGATGATGCA
CCTGACCCGGTTTCACTCTCA
ACGACTCACTCCACTCAGAGG
CCACCGCGCTGTGCACGTCC
ACCAAGATOUTTACCACACTTA
CACATCACTCTCAGAGCTCAC
TCCACTCAGAGGCCACCCGCG
GCTGTGCACTCCACACGATC
CTTACCACCTTACACATTACC
ATATATCCACTTACCACACATA
CCTTACCAGATATCCACTACC
ACCATATATCCACTTACACAC
ACCTATTATTATTACCGAGGGA
GAGGGGTGACCACACTGTGACA
```


Vertebrate Genome Sequences




Mouse




Rat




Chicken




Chimpanzee




Dog




Macaque




Monodelphis




Orangutan




Marmoset




Horse




Cow




Platypus



Xenopus

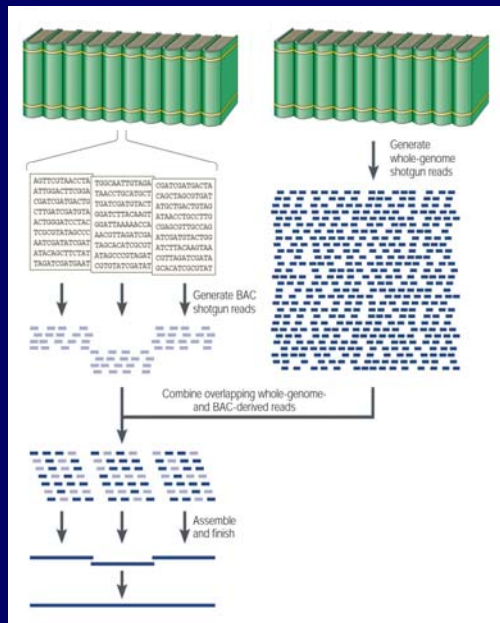


Zebrafish



Pufferfish

Hybrid Shotgun Sequencing



Green (2001)

Diverse Landscape of Genome Sequencing

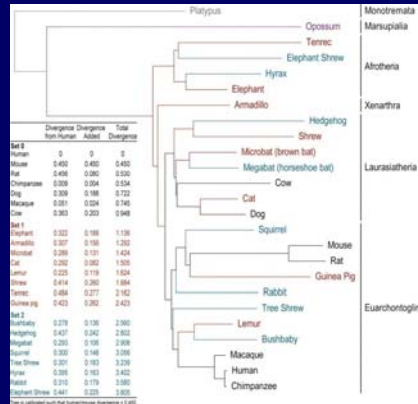
Human	=====
Mouse	=====
Rat	=====
Pufferfish	=====
Zebrafish	=====
Chicken	=====
Chimpanzee	=====
Dog	=====
Cow	=====
Xenopus	=====
Monodelphis	=====
Macaque	=====
Platypus	=====
Marmoset	=====
etc....	=====

Low-Redundancy, Whole-Genome Shotgun Sequencing

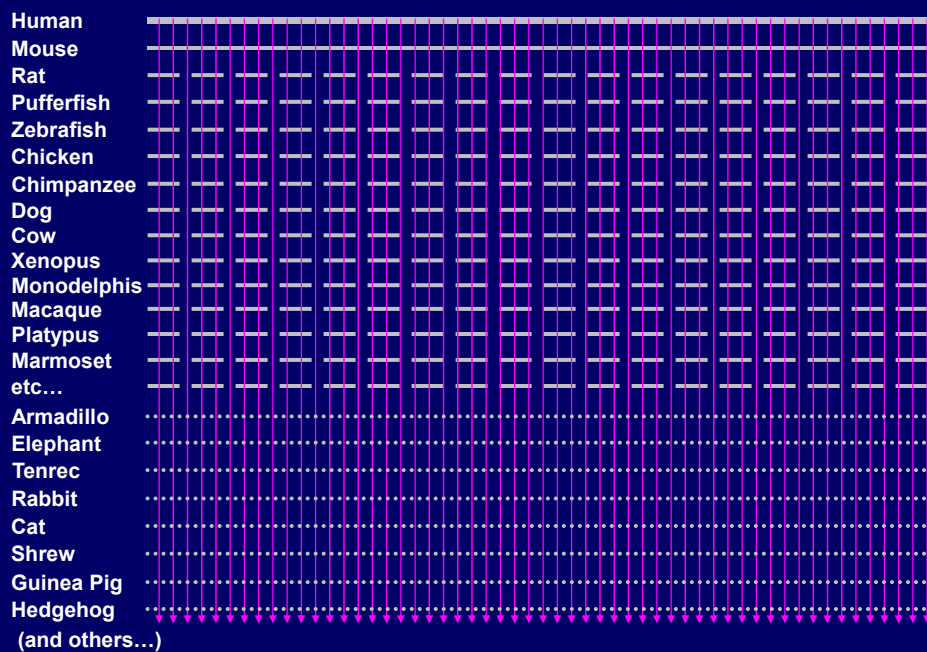
An initial strategy for the systematic identification of functional elements in the human genome by low-redundancy comparative sequencing

Elliott H. Margulies^{*†}, Jade Vinson^{†‡}, NISC Comparative Sequencing Program^{*§}, Webb Miller[¶], David B. Jaffe[¶], Kerstin Lindblad-Toh[¶], Jean Chang[¶], Eric D. Green^{*¶}, Eric S. Lander[¶], James C. Mullikin^{*¶¶}, and Michele Clamp^{*¶¶}

Margulies et al. (2005)



Diverse Landscape of Genome Sequencing



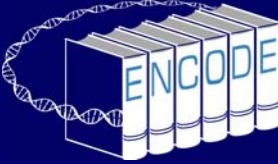
Multi-Species Sequence Comparisons



Multi-Species Conserved Sequences (MCSs)

Margulies et al. (2003)
Thomas et al. (2003)

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

SPECIAL SECTION GENES IN ACTION
VIEWPOINT

The ENCODE (ENCyclopedia Of DNA Elements) Project

The ENCODE Project Consortium*†

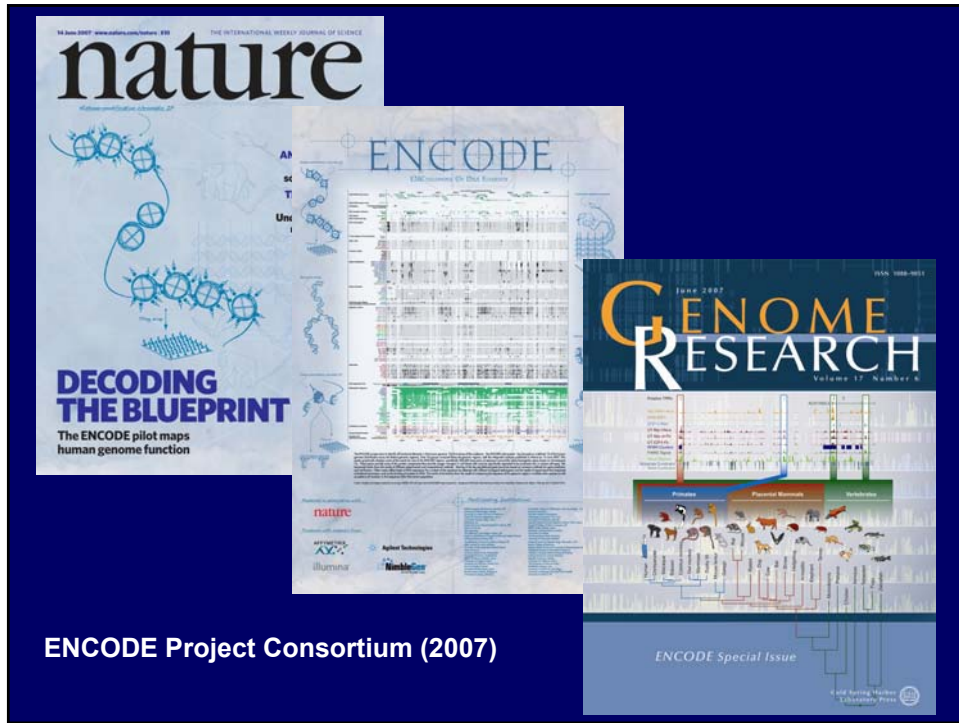
ENCODE Project Consortium (2004)

The diagram illustrates the ENCODE project's approach to identifying DNA elements. It shows a DNA strand with various elements: Long-range regulatory elements (enhancers, repressors/silencers, insulators), cis-regulatory elements (promoters, transcription factor binding sites), and a Gene. The workflow involves DNase Digestion, Reporter Assays, ChIP-chip, Microarray Hybridization, and Computational Predictions and RT-PCR. Epigenetic modifications (CH₃CO, CH₃) and DNA Replication Sites are also indicated.

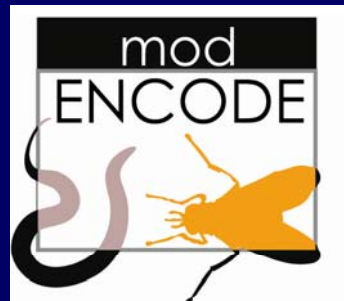
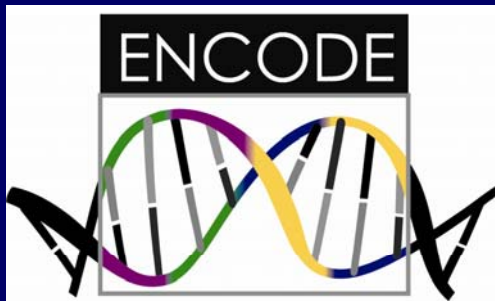
ENCODE Project: Web Sites

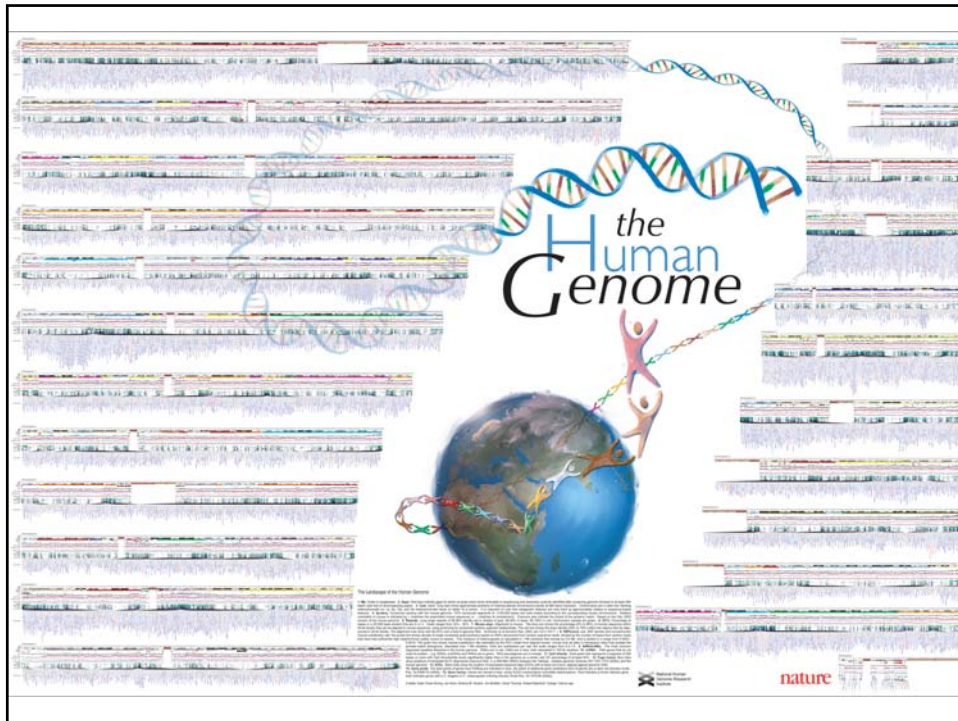
The left screenshot shows the ENCODE Target Regions (January 2004) page, which provides background information and lists identified random picks. The right screenshot shows the UCSC Genome Browser on Human July 2003 Freeze, displaying a genomic track for a specific region.

genome.gov/ENCODE genome.ucsc.edu/ENCODE

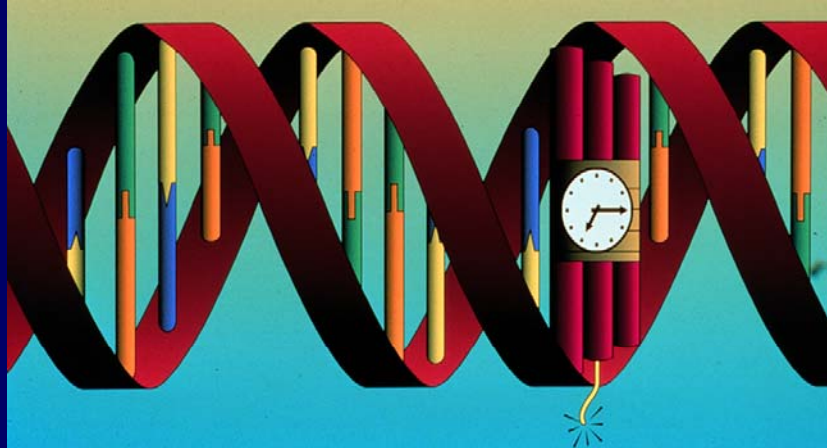


Expanding ENCODE Portfolio

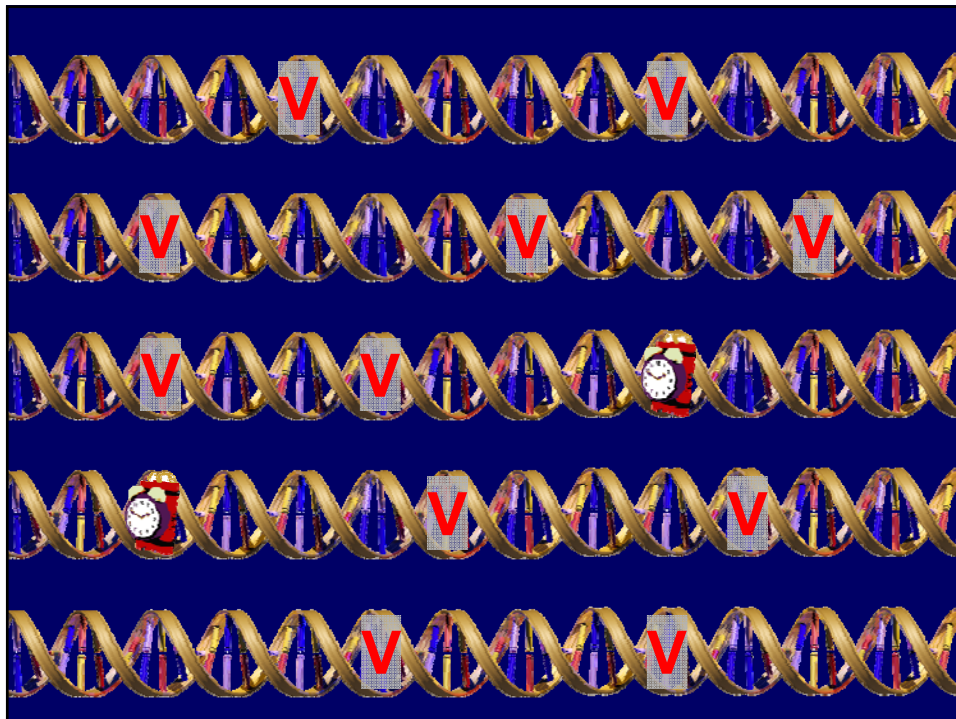


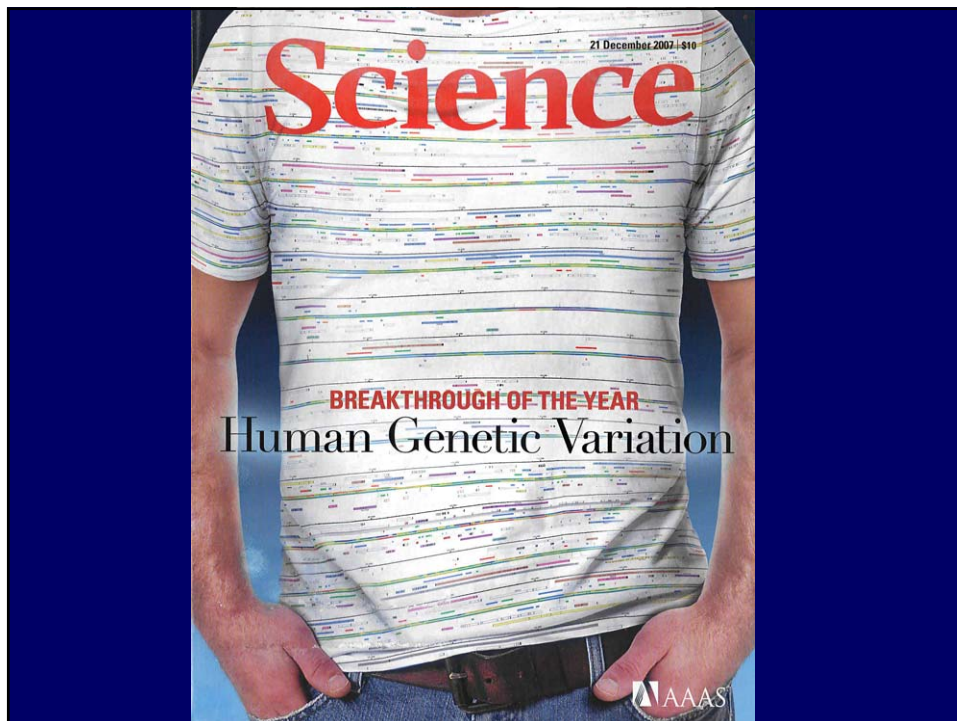


All humans are ~99.9% identical at the DNA sequence level, and yet...



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







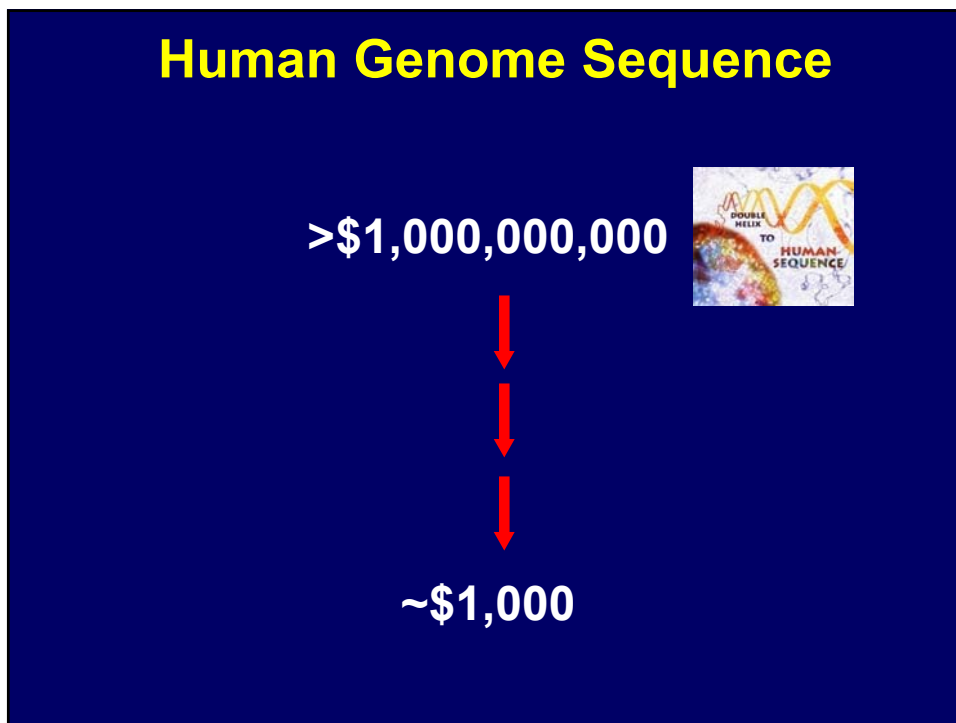
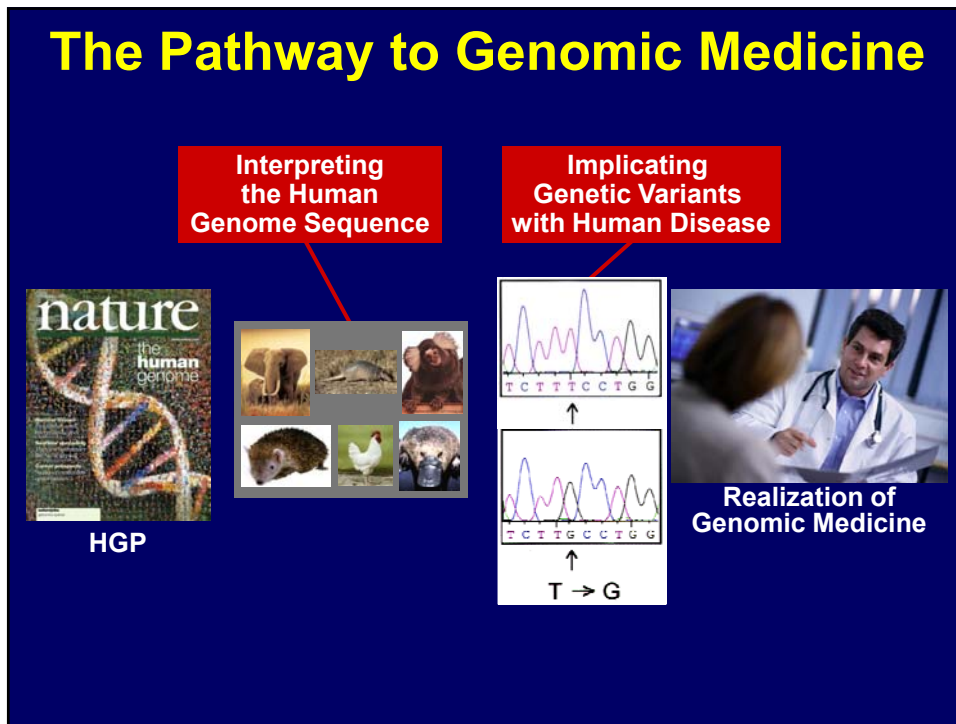
The Pathway to Genomic Medicine

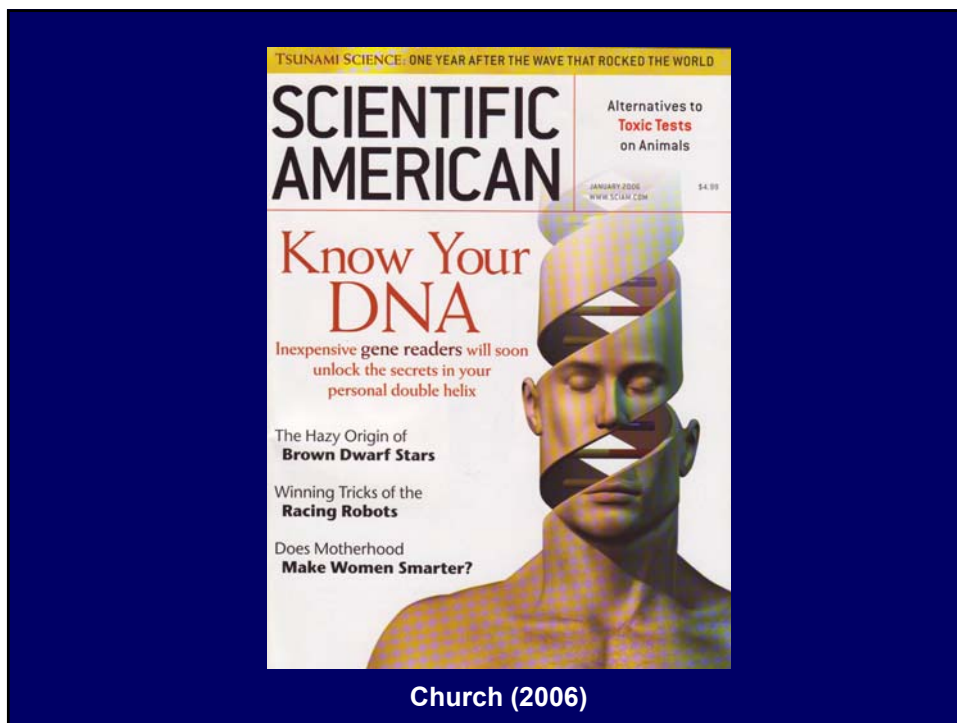
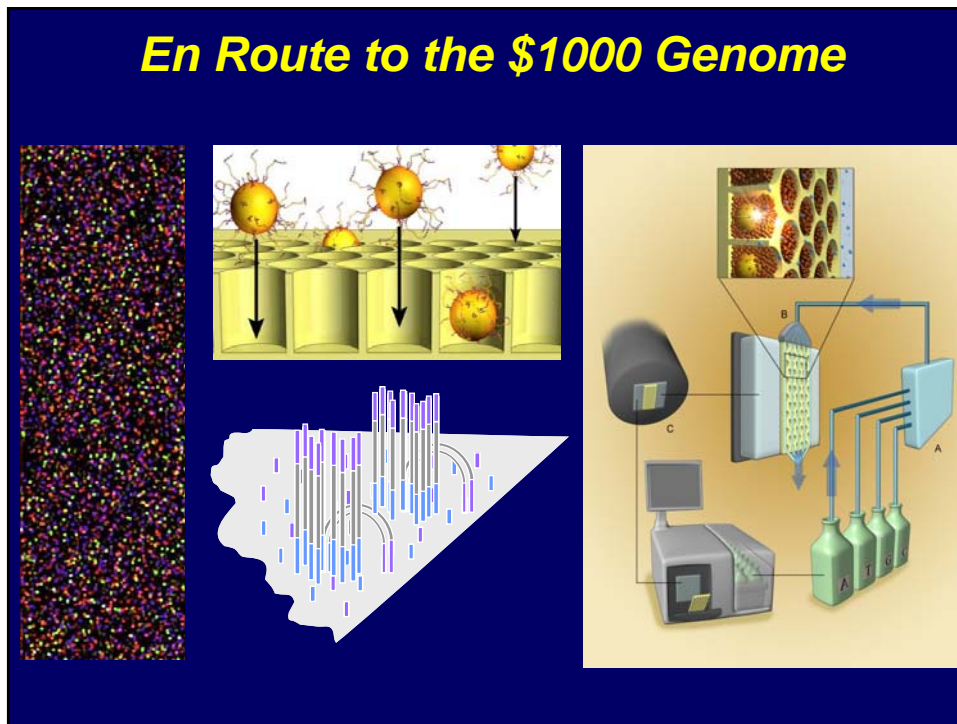


HGP



Realization of Genomic Medicine

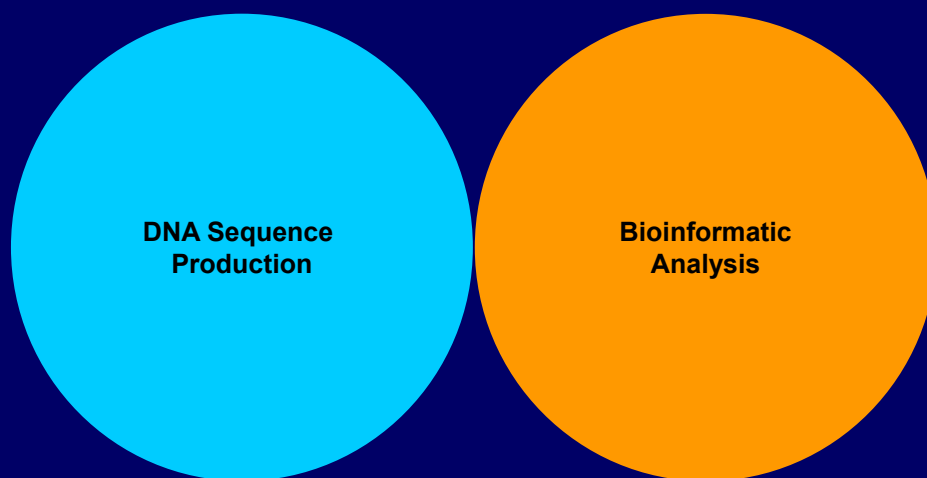




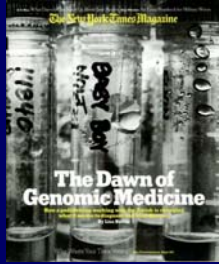
Realities of New DNA Sequencing Technologies...



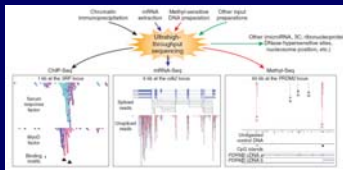
Changing Infrastructure Requirements



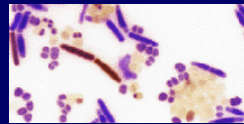
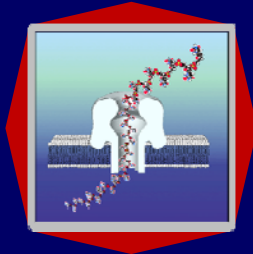
Expanding Universe of Sequence-Based Explorations



Collins and Barker (2007)



Wold and Myers (2008)



Turnbaugh et al. (2007)



Enard and Paabo (2004)



Green et al. (2006)

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.
- 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.
- Collins FS and Barker AD (2007). Mapping the cancer genome: pinpointing the genes involved in cancer will help chart a new course across the complex landscape of human malignancies. *Sci Am* 296:50-57.
- Enard W and Paabo S (2004). Comparative primate genomics. *Annu Rev Genomics Hum Genet* 5:351-378.
- ENCODE Project Consortium (2004). The ENCODE (ENCyclopedia Of DNA Elements) Project. *Science* 306:636-640.
- ENCODE Project Consortium (2007). Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project. *Nature* 447:799-816.
- 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.
- Green RE et al. (2006). Analysis of one million base pairs of Neanderthal DNA. *Nature* 444:330-336.
- Hillier LW et al. (2004). Sequence and comparative analysis of the chicken genome provide unique perspectives on vertebrate evolution. *Nature* 432:695-716.
- International HapMap Consortium (2005). A haplotype map of the human genome. *Nature* 437:1299-1320.

- International HapMap Consortium (2007). A second generation human haplotype map of over 3.1 million SNPs. *Nature* 449:851-861.
- 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 EH et al. (2003). Identification and characterization of multi-species conserved sequences. *Genome Res* 13:2507-2518.
- Margulies EH 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.
- 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.
- Mikkelsen TS et al. (2007). Genome of the marsupial *Monodelphis domestica* reveals innovation in non-coding sequences. *Nature* 447:167-177.
- 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.
- Rhesus Macaque Genome Sequencing and Analysis Consortium (2007). Evolutionary and biomedical insights from the rhesus macaque genome. *Science* 316:222-234.
- Thomas JW et al. (2003). Comparative analyses of multi-species sequences from targeted genomic regions. *Nature* 424:788-793.
- Turnbaugh PJ et al. (2007). The human microbiome project. *Nature* 449:804-810.
- 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.
- Wold B and Myers RM (2008). Sequence census methods for functional genomics. *Nat Methods* 5:19-21.