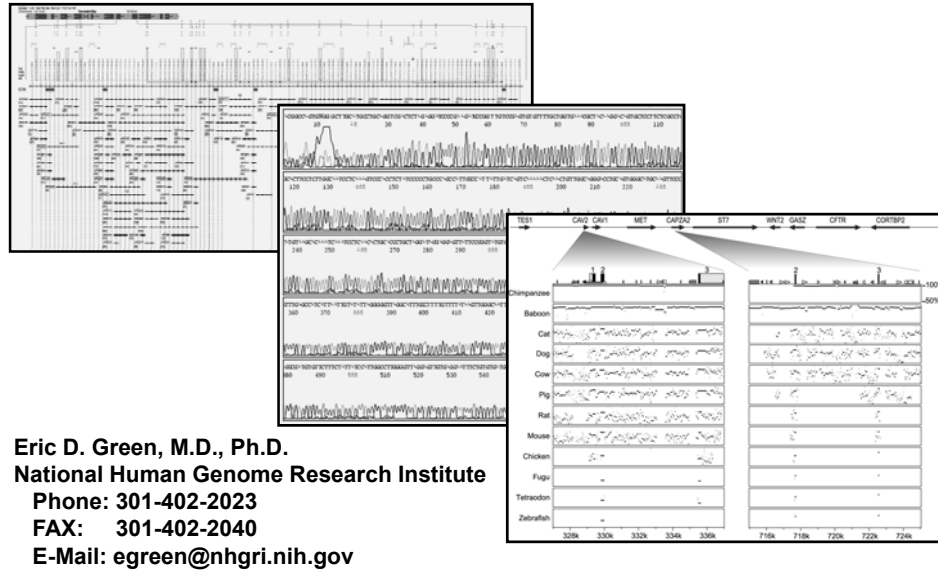


## Techniques for Analyzing Genomes I



## Foundational Milestones in Genetics & Genomics



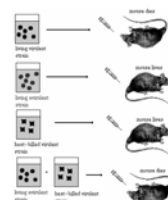
**Mendel**

**1865**



**Miescher**

**1871**



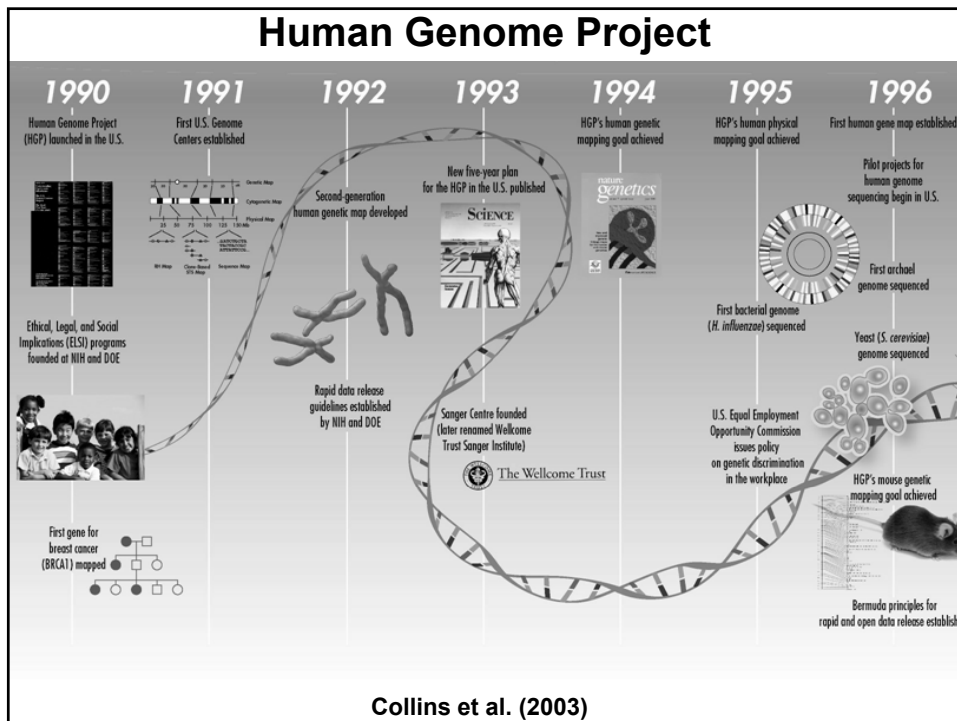
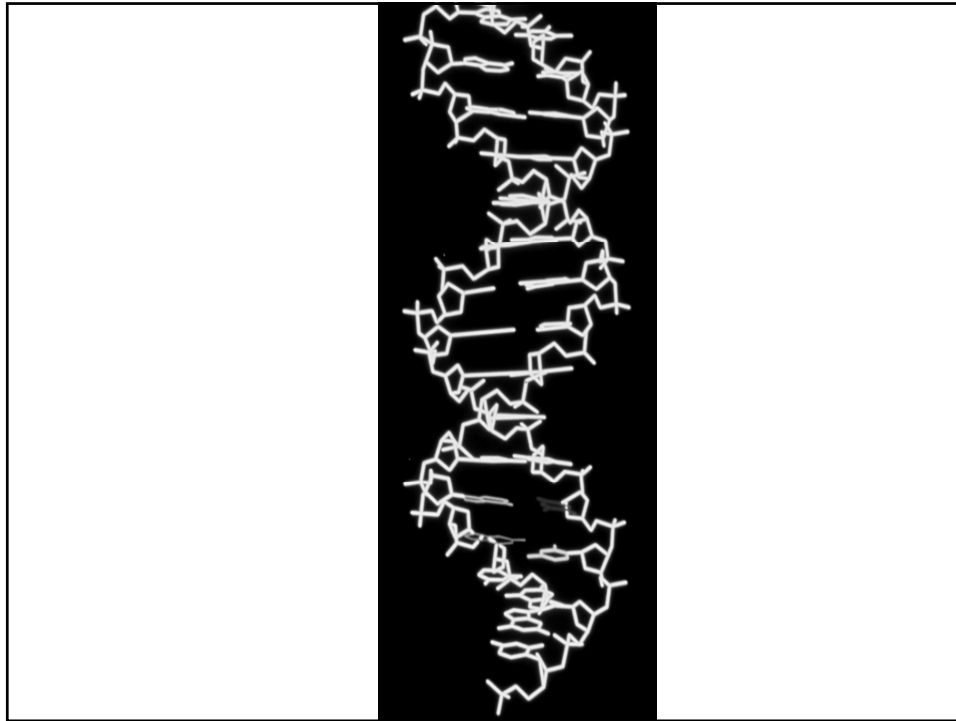
**Avery**

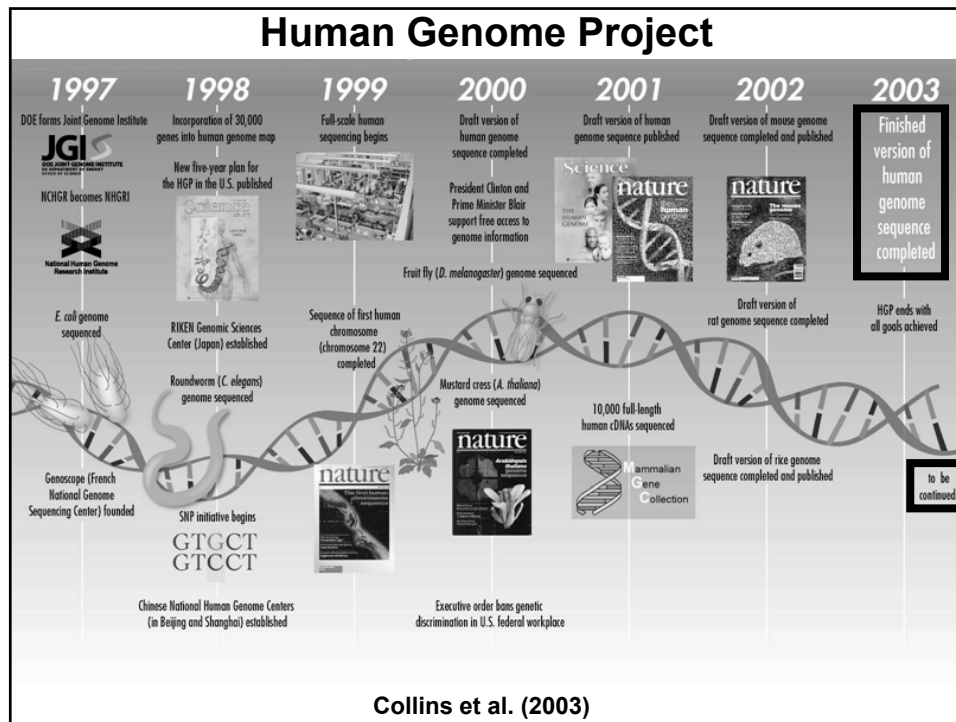
**1944**



**Watson & Crick**

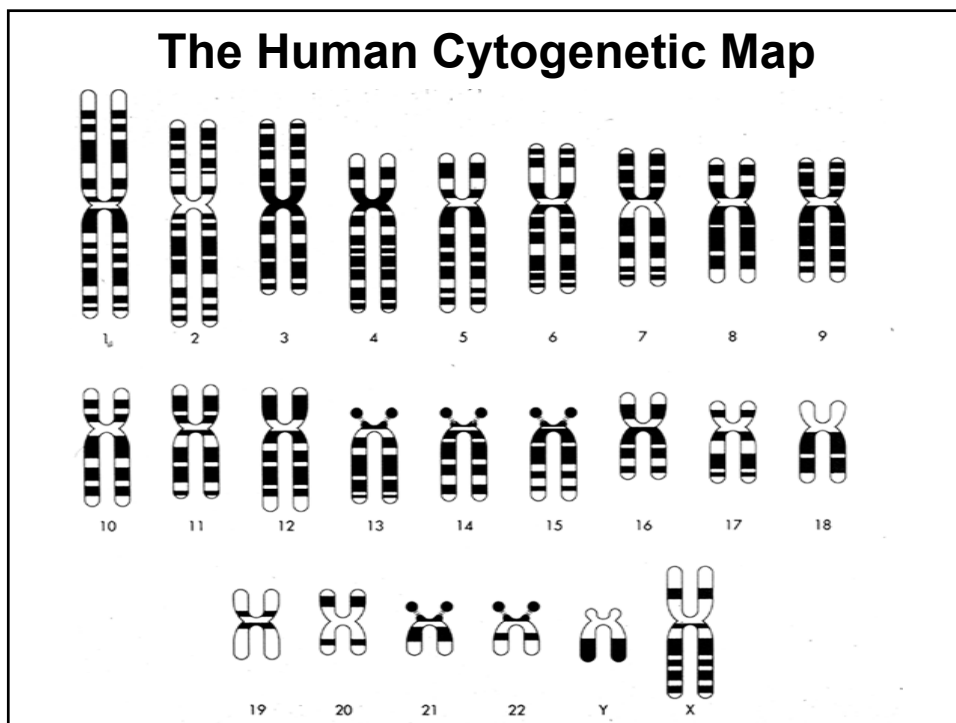
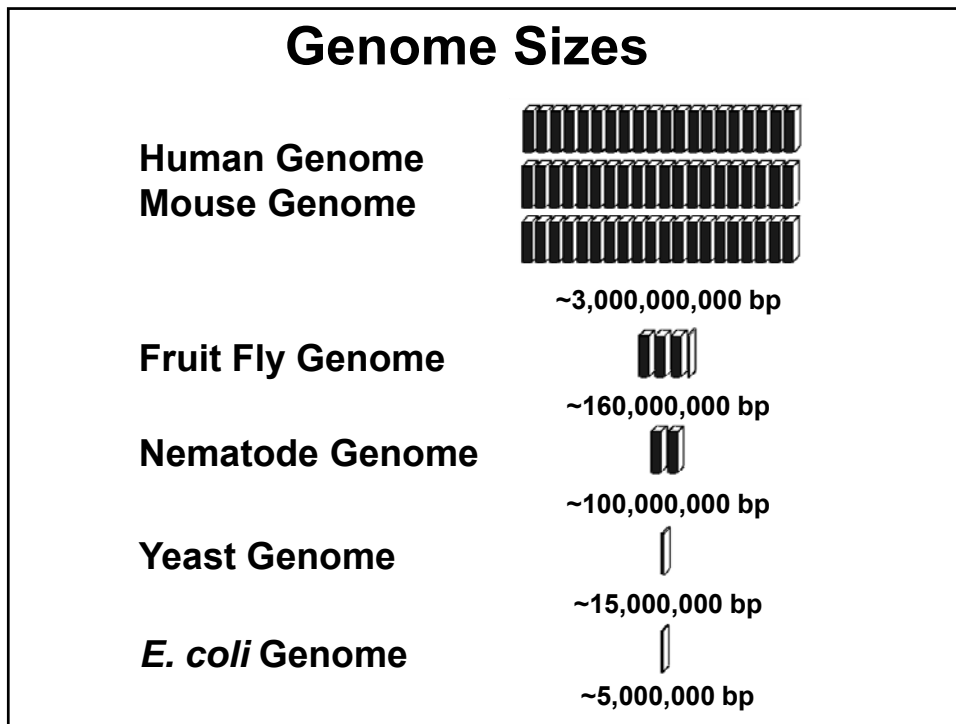
**1953**

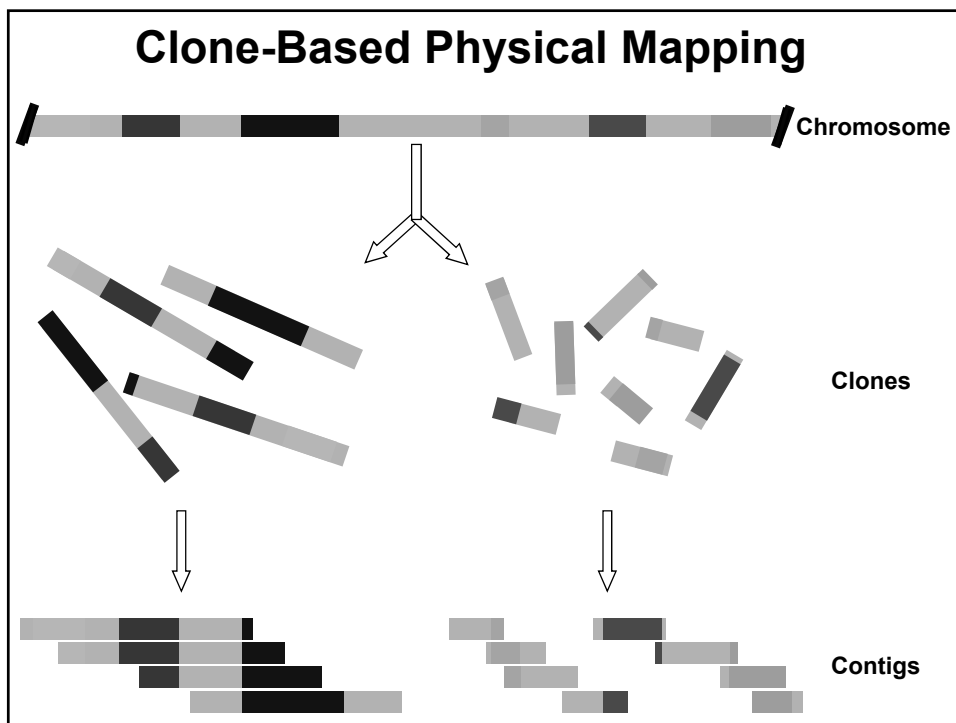
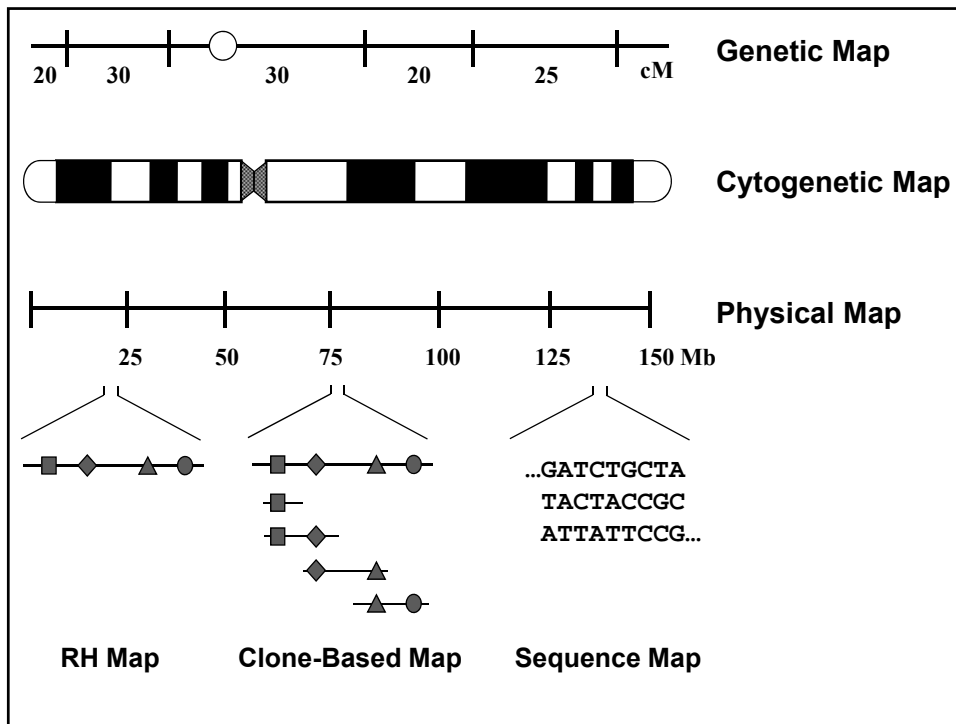


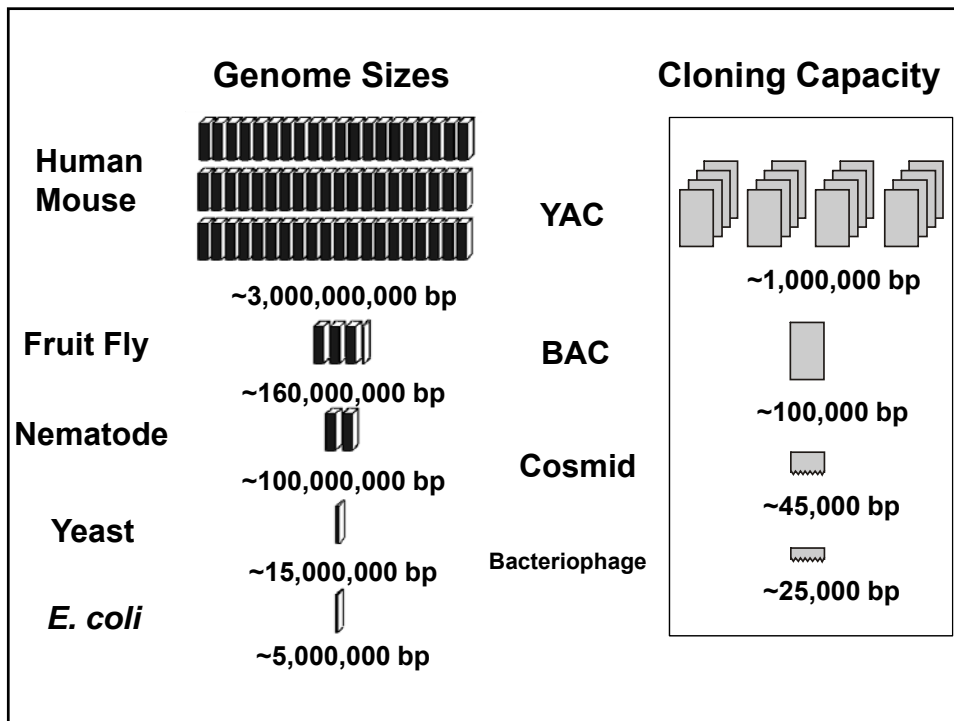
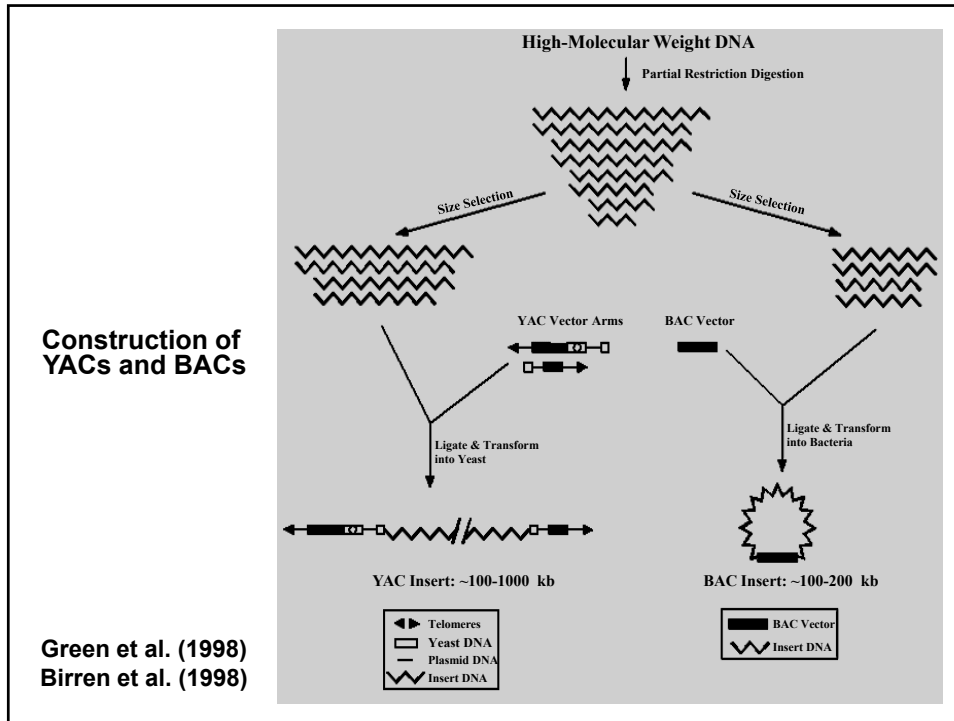


## Outline

- 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

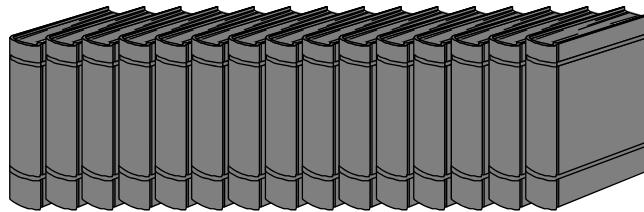




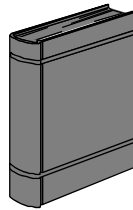


## Bacterial Artificial Chromosomes (BACs)

- 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](http://www.chori.org/bacpac))



**Genome**  
(~3000 Mb)



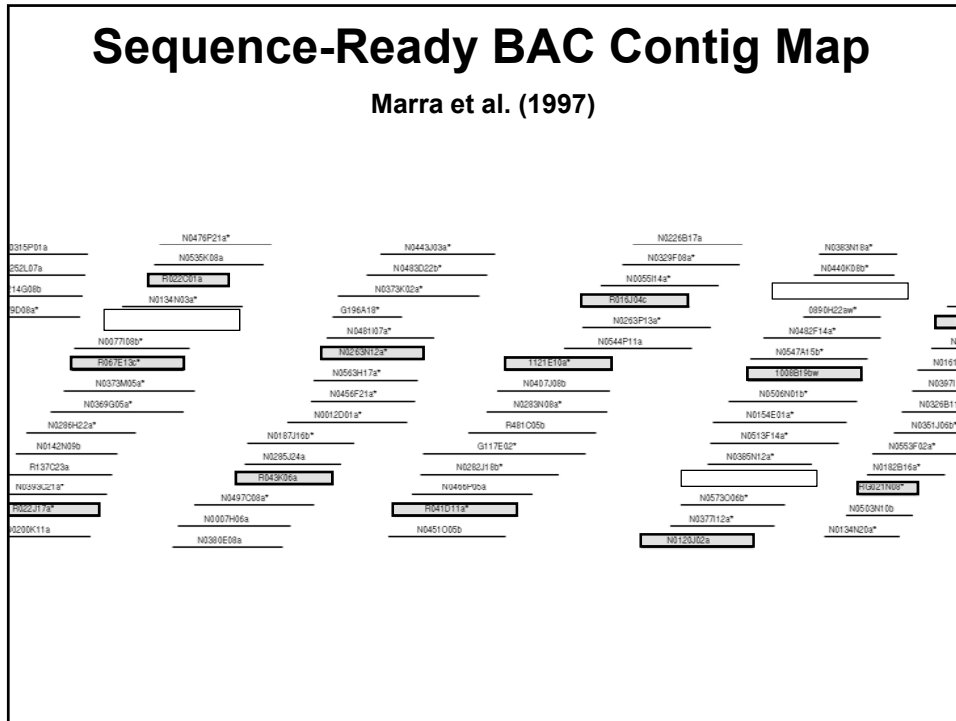
**Chromosome**  
(~130 Mb)

GI	GI	GI	GI	GI	GATCCTCTAGAATCTC
GI	GI	GI	GI	GI	GAGATCTCTGAGAGTC
GT	GT	GT	GT	GT	GTGGGAACTGTTTGA
TT	TT	TT	TT	TT	TGTGACTAGCCACAGT
TT	TT	TT	TT	TT	TGTGACTAGCCACAGT
TT	TT	TT	TT	TT	TACTGTGTGAGAGATGT
AT	AT	AT	AT	AT	ATGATGGACCTGACCC
GI	GI	GI	GI	GI	GGGTTTCACTCTCAAC
GI	GI	GI	GI	GI	GACTCACTCCACTTCA
CT	CT	CT	CT	CT	CCGGTTAGATACAG
GI	GI	GI	GI	GI	GAGGCCCAACCACCCT
GT	GT	GT	GT	GT	GTGCACGTCACACC

**YAC**  
(~0.5-1.0 Mb)

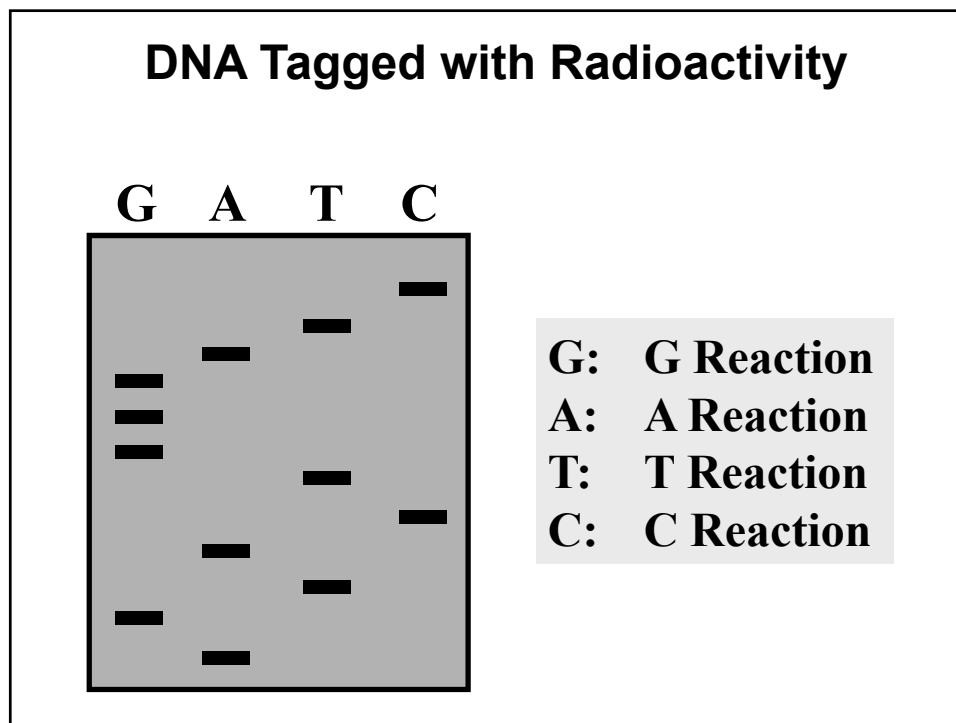
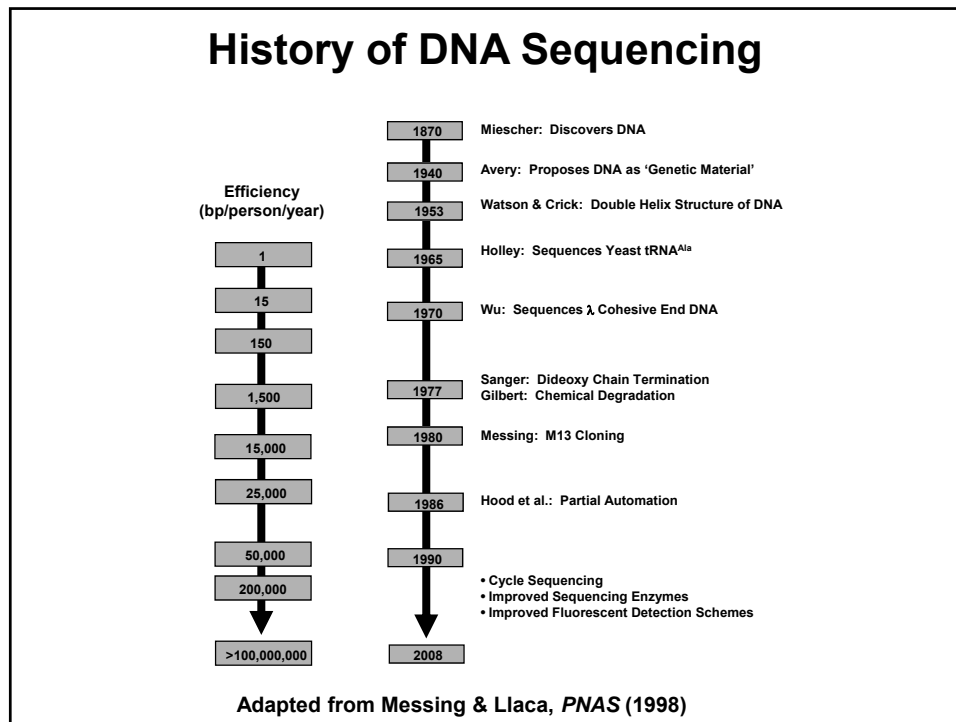
GATCCTCTAGAATCTC
GAGATCTCTGAGAGTC
GTGGGAACTGTTTGA
TGTGACTAGCCACAGT
TAGGTATTGGGCAATT
TACTGTGTGAGAGATGT
ATGATGGACCTGACCC
GGGTTTCACTCTCAAC
GACTCACTCCACTTCA
CCGGTTAGATACAG
GAGGCCCAACCACCCT
GTGCACGTCACACC

**BAC**  
(~0.1-0.2 Mb)



# DNA Sequencing



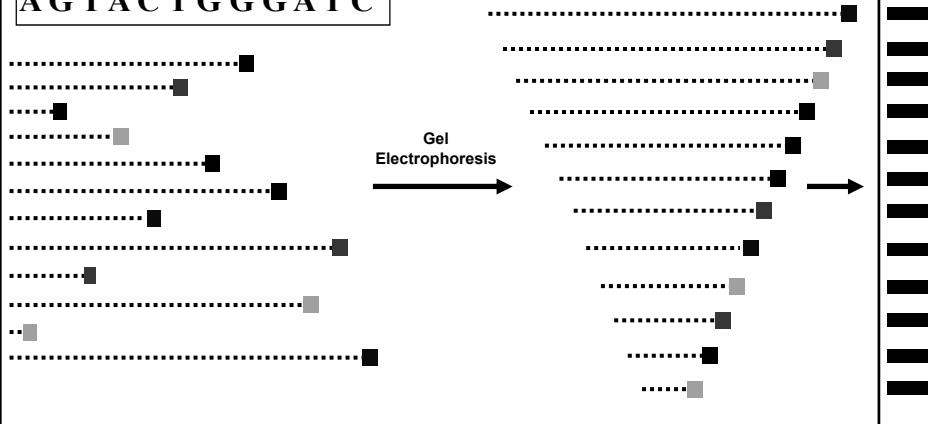


## Radioactive Sequencing

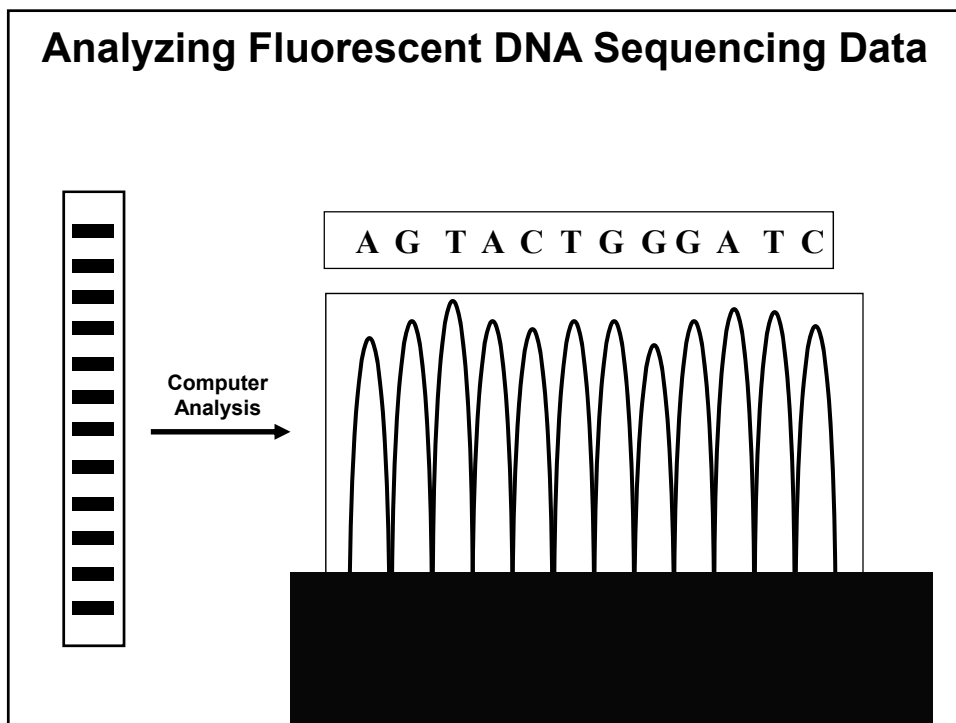
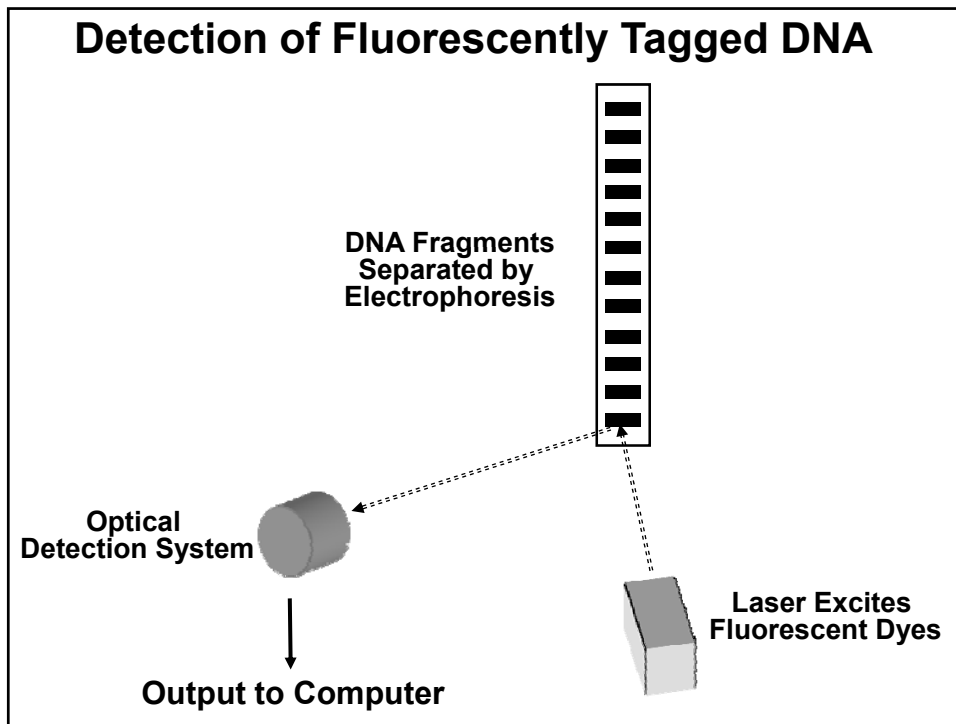


## Fluorescent DNA Sequencing

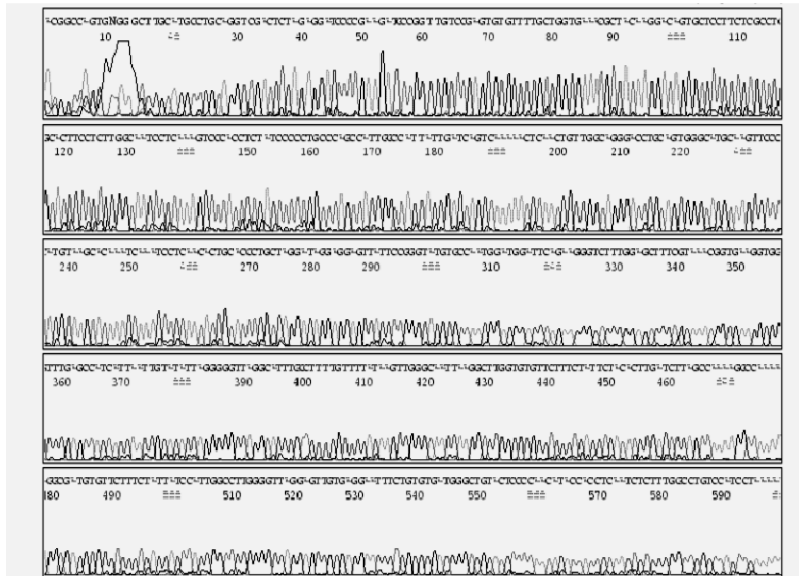
AGTACTGGGATC



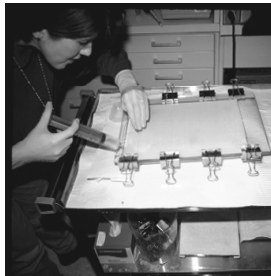
Wilson & Mardis (1997)



## Fluorescent DNA Sequencing Results



## Slab Gel-Based DNA Sequencing Instruments

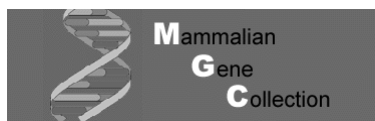


## Capillary-Based DNA Sequencing Instruments



## Large-Scale cDNA Sequencing

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



[mgc.nci.nih.gov](http://mgc.nci.nih.gov)

Gerhard et al. (2004)

# Large-Scale Genome Sequencing



# Shotgun Sequencing

Wilson & Mardis (1997)  
Green (2001)

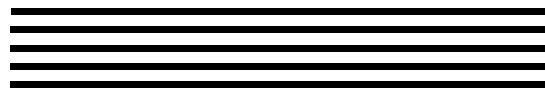
## Subclone Construction

```
GATCCTTAGAATCTC
GAGATCTTAGAATCTC
GTGGAAACTGTGTGA
TTTGTACTACCAAGT
TACCTGTAGAGATGT
ATGATGCACCTTGACC
GGATTTGATCTTAGAC
GACTCACTCCACTCA
GAGGCCACCCCTGCT
GTGCACCTCCACCAC
GATTTATACATTTTA
AATCTTAGAATGAGA
```

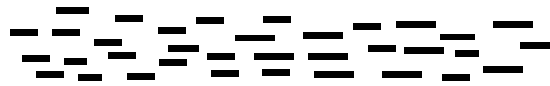
————— BAC DNA



Prepare Multiple Copies



Randomly Fragment

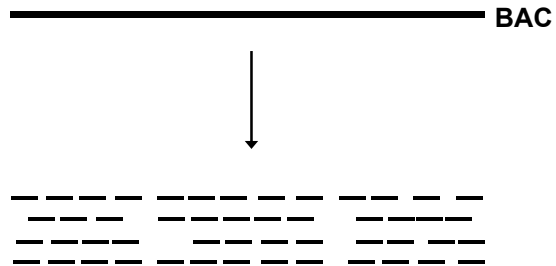


Subclone Fragments

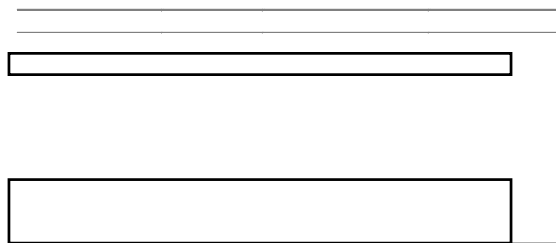


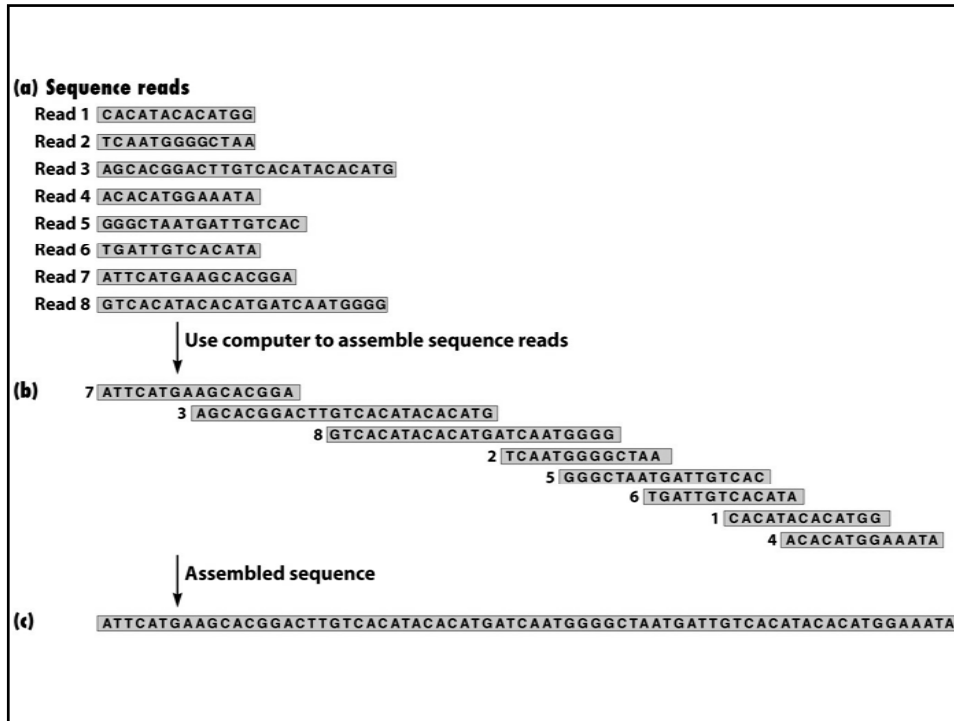
```
GA GA GA GATCCTTAGAATCTC
GA GA GA GAGATCTTAGAATCTC
GA GA GA GTGGAAACTGTGTGA
TT TT TT TTTGTACTACCAAGT
AT AT AT ATGATGCACCTTGACC
GA GA GA GGATTTGATCTTAGAC
GA GA GA GACTCACTCCACTCA
GA GA GA GAGGCCACCCCTGCT
GA GA GA GTGCACCTCCACCAC
GA GA GA GATTTATACATTTTA
AT AT AT AATCTTAGAATGAGA
```

## Shotgun Sequencing Strategy



## Poisson Calculations

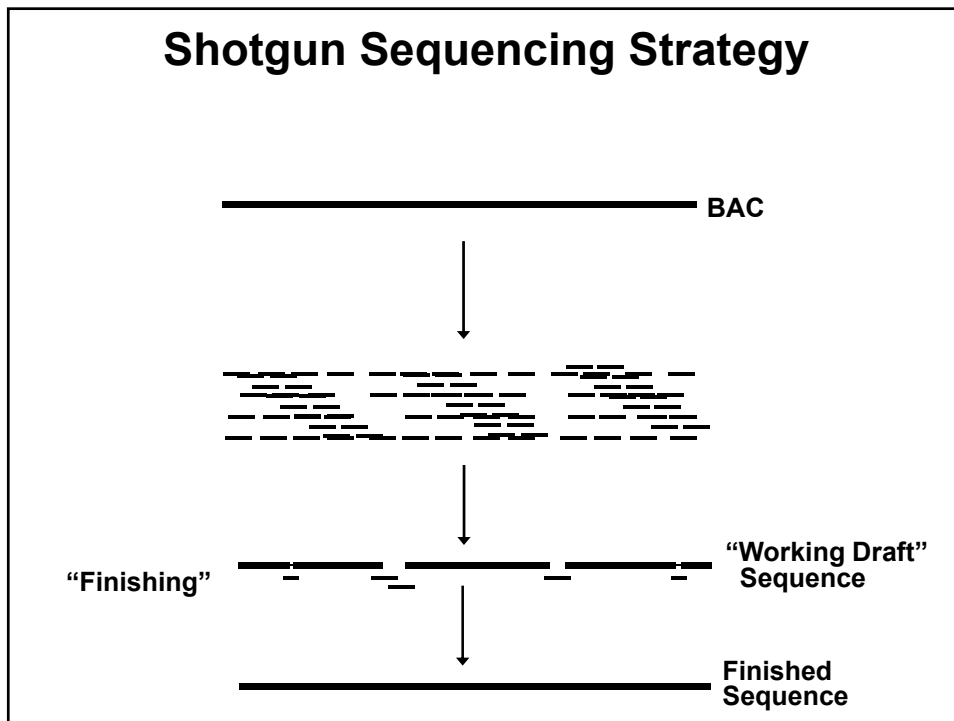
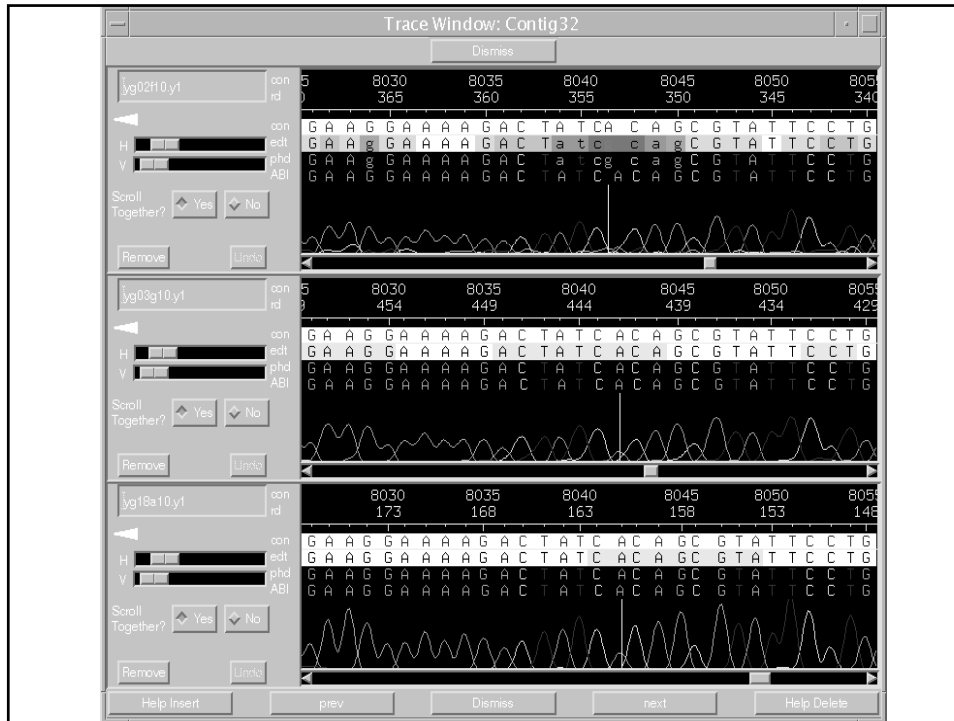




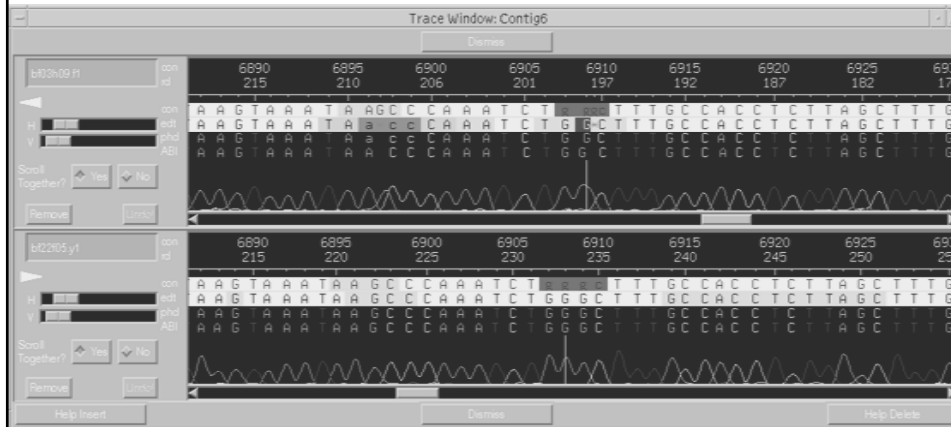
## Shotgun Sequence Assembly

“Consed” (Gordon et al., 1998)



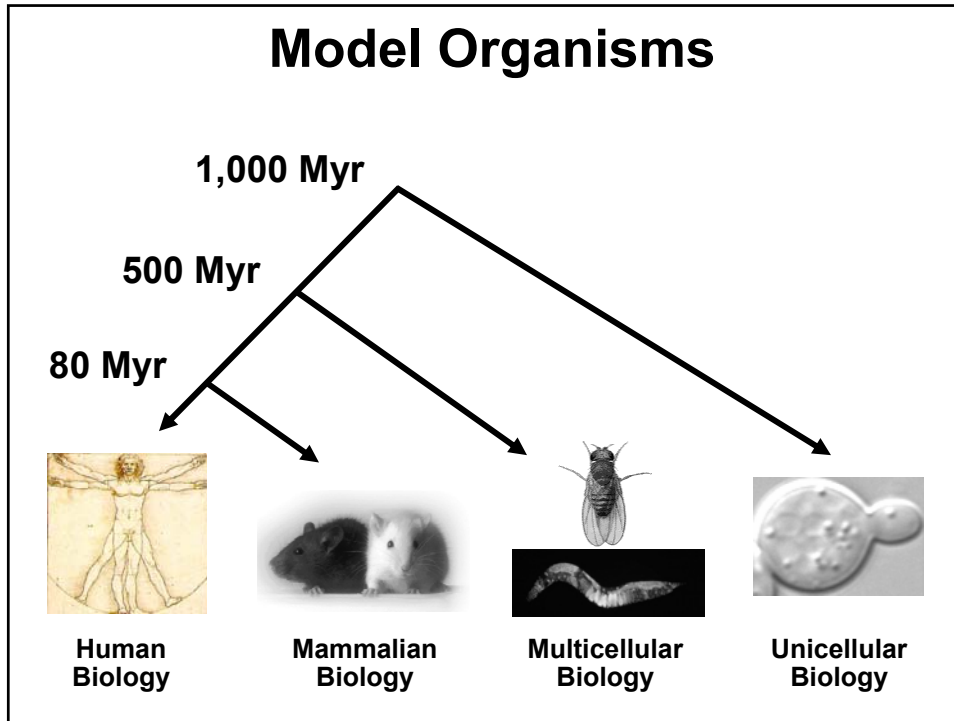


## Sequence Finishing: Resolving Ambiguities





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

## Historically Significant Genome Sequencing Projects



## Microbial Genome Sequences

Comprehensive Microbial Resource

Home Genome Tools Searches Comparative Tools Lists Downloads Carts

Search Locus for  Go

---

**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**

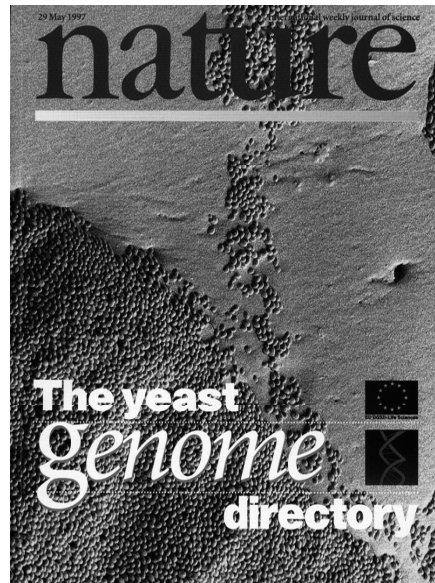


**March 13, 2007: CAMERA** is a web resource for metagenomic research. CAMERA's debut coincides with the publication of the [Global Ocean Sampling](#) expedition's extensive dataset cataloging over 6 million new genes from uncultured marine microbes. Come visit [CAMERA](#), and see our growing collection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

**Latest Releases**  
Data Release: [21.0](#)

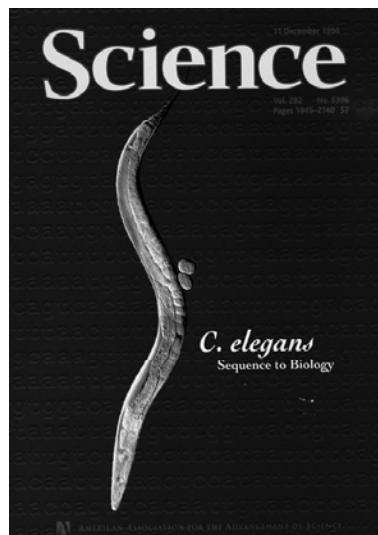
[www.tigr.org](http://www.tigr.org)

## First Eukaryotic Genome Sequence



Goffeau et al. (1997)

## First Animal Genome Sequence

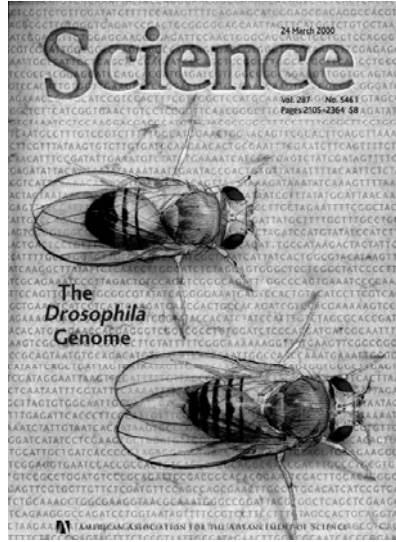


Genome Sequence of the Nematode *C. elegans*:  
A Platform for Investigating Biology

The *C. elegans* Sequencing Consortium\*

*C. elegans* Sequencing Consortium (1998)

# Second Animal Genome Sequence



THE DROSOPHILA GENOME  
REVIEW

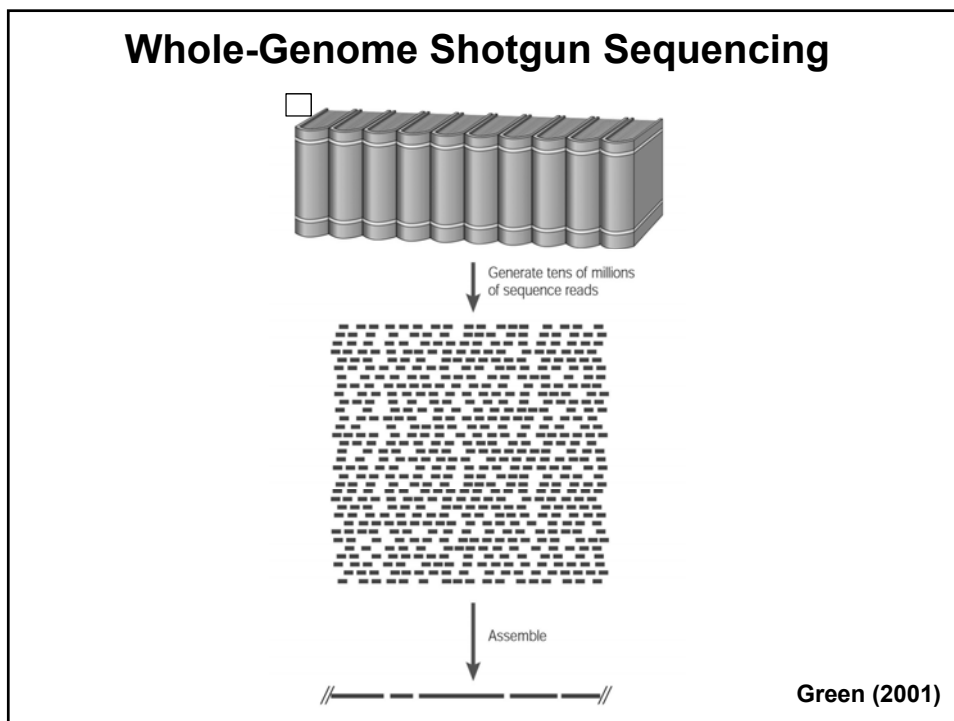
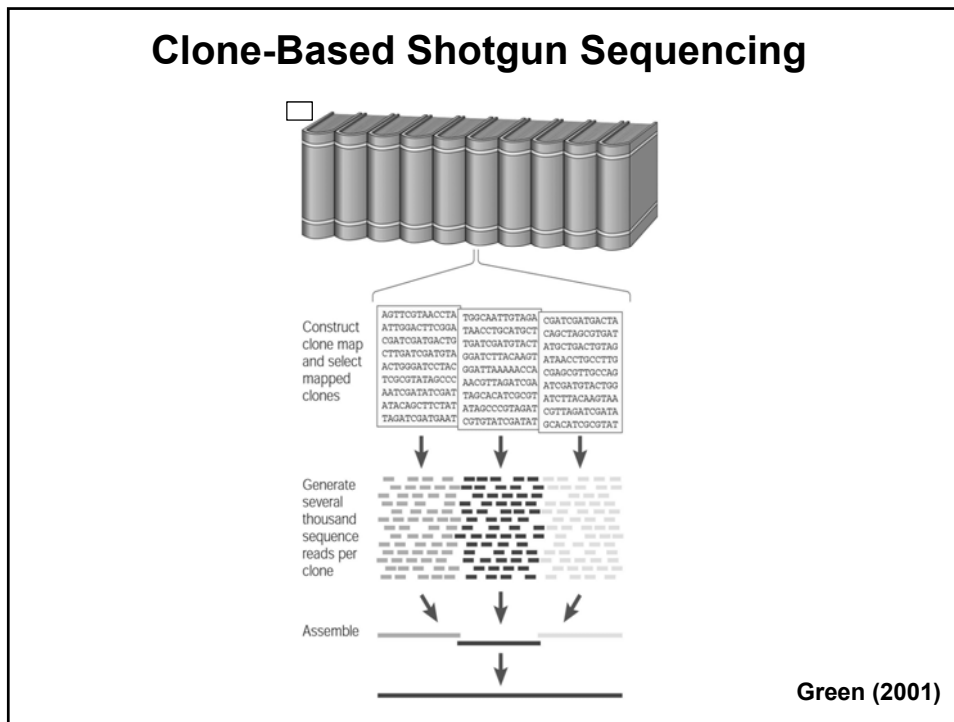
**The Genome Sequence of *Drosophila melanogaster***

Mark D. Adams,<sup>1\*</sup> Susan E. Celniker,<sup>2</sup> Robert A. Holt,<sup>1</sup> Cheryl A. Evans,<sup>1</sup> Jeannine D. Gockyne,<sup>1</sup> Peter G. Amanatides,<sup>1</sup> Steven E. Scherer,<sup>3</sup> Peter W. Li,<sup>1</sup> Roger A. Hoskins,<sup>2</sup> Richard F. Galie,<sup>1</sup> Reed A. George,<sup>4</sup> Suzanna E. Lewis,<sup>1</sup> Stephen Richards,<sup>5</sup> Michael Ashburner,<sup>6</sup> Scott N. Henderson,<sup>7</sup> Granger G. Sutton,<sup>8</sup> Jennifer R. Wortman,<sup>1</sup> Mark D. Vandell,<sup>1</sup> Qing Zhang,<sup>1</sup> Lin X. Chen,<sup>1</sup> Rhonda C. Brandon,<sup>1</sup> Yu-Hui C. Rogers,<sup>1</sup> Robert G. Blasziak,<sup>1</sup> Mark Champe,<sup>9</sup> Berret D. Pfeiffer,<sup>1</sup> Kenneth H. Wan,<sup>1</sup> Clare Doyle,<sup>1</sup> Evan G. Baxter,<sup>1</sup> Gregg Helt,<sup>1</sup> Catherine R. Nelson,<sup>1</sup> George L. Gabor Miklos,<sup>1</sup> Josep F. Abril,<sup>10</sup> Anna Aghayani,<sup>11</sup> Hui-Jin An,<sup>1</sup> Cynthia Andrawe-Pfaendler,<sup>12</sup> Dorothea Baldwin,<sup>1</sup> Richard M. Balow,<sup>1</sup> Anand Basu,<sup>1</sup> James Baxterdale,<sup>1</sup> Leyla Bayraktaroglu,<sup>13</sup> Ellen M. Beasley,<sup>1</sup> Karen Y. Beeson,<sup>1</sup> P. V. Benos,<sup>14</sup> Benjamin P. Berman,<sup>15</sup> Deepali Bhandari,<sup>16</sup> Slava Bolshakov,<sup>17</sup> Dana Borkova,<sup>18</sup> Michael R. Botchan,<sup>19</sup> John Bouck,<sup>20</sup> Peter Brokstein,<sup>21</sup> Philippe Brotier,<sup>22</sup> Kenneth C. Burtis,<sup>1</sup> Dana A. Busam,<sup>1</sup> Heather Butler,<sup>23</sup> Edouard Cadieu,<sup>24</sup> Angela Center,<sup>25</sup> Ishwar Chandra,<sup>1</sup> J. Michael Cherry,<sup>18</sup> Simon Cawley,<sup>18</sup> Carl Dahlke,<sup>1</sup> Lionel B. Davenport,<sup>1</sup> Peter Davies,<sup>1</sup> Beatriz de Pablos,<sup>26</sup> Arthur Delcher,<sup>27</sup> Zuoming Deng,<sup>28</sup> Anne Deslattes Mays,<sup>1</sup> Ian Dew,<sup>1</sup> Suzanne M. Dietz,<sup>1</sup> Kristina Dodson,<sup>1</sup> Lisa E. Doup,<sup>1</sup> Michael Downes,<sup>21</sup> Shannon Dugan-Rocha,<sup>29</sup> Boris C. Dunkov,<sup>30</sup> Patrick Dunn,<sup>31</sup> Kenneth J. Durbin,<sup>32</sup> Carlos C. Evangelista,<sup>33</sup> Concepcion Ferraz,<sup>34</sup> Steven Ferrieri,<sup>1</sup> Wolfgang Fleischmann,<sup>1</sup> Carl Foster,<sup>1</sup> Andrei E. Gabrielian,<sup>1</sup> Neha S. Garg,<sup>1</sup> William M. Gelbart,<sup>35</sup> Ken Glasser,<sup>1</sup> Anna Glöckl,<sup>36</sup> Fangcheng Gong,<sup>1</sup> J. Harley Gorrell,<sup>37</sup> Zhiping Gu,<sup>38</sup> Fing Guan,<sup>39</sup> Michael Harris,<sup>40</sup> Naomi L. Harris,<sup>41</sup> Damon Harvey,<sup>42</sup> Thomas J. Heiman,<sup>43</sup> Judith R. Hernandez,<sup>44</sup> Jarrett Houck,<sup>45</sup> Damon Houston,<sup>46</sup> Kathryn A. Houston,<sup>47</sup> Timothy J. Howland,<sup>48</sup> Hing-Hui Wei,<sup>49</sup> Chinyere Ibegwam,<sup>50</sup> Mensa Jalali,<sup>51</sup> Francis Kalish,<sup>52</sup> Gary H. Karpen,<sup>53</sup> Zhaod Ke,<sup>54</sup> James A. Kennison,<sup>55</sup> Karen A. Ketchum,<sup>56</sup> Bruce E. Kimmel,<sup>57</sup> Chinnappa D. Kodira,<sup>58</sup> Cheryl Kraft,<sup>59</sup> Saul Kravitz,<sup>60</sup> David Kulp,<sup>61</sup> Zhongyu Lei,<sup>62</sup> Paul Lasko,<sup>63</sup> Yiding Lei,<sup>64</sup> Alexander A. Levitsky,<sup>65</sup> Jayin Li,<sup>66</sup> Zhenyu Li,<sup>67</sup> Yong Liang,<sup>68</sup> Xiaoying Lin,<sup>69</sup> Xiangjun Liu,<sup>70</sup> Bettina Hates,<sup>71</sup> Tina C. McIntosh,<sup>72</sup> Michael P. McLeod,<sup>73</sup> Duncan McPherson,<sup>74</sup> Genady Merkulov,<sup>75</sup> Natalia V. Milshina,<sup>76</sup> Clark Moberly,<sup>77</sup> Joe Morris,<sup>78</sup> Ali Moshrefi,<sup>79</sup> Stephen M. Mount,<sup>80</sup> Mei Mo,<sup>81</sup> Brian Murphy,<sup>82</sup> Lee Murphy,<sup>83</sup> Donna M. Muzny,<sup>84</sup> David L. Nelson,<sup>85</sup> David R. Nelson,<sup>86</sup> Keith A. Nelson,<sup>87</sup> Katherine Nixon,<sup>88</sup> Deborah R. Nusbaum,<sup>89</sup> Joanne M. Pacifico,<sup>90</sup> Michael Palazzolo,<sup>91</sup> Gijung S. Pitman,<sup>92</sup> Sun Pan,<sup>93</sup> John Pollard,<sup>94</sup> Vinita Puri,<sup>95</sup> Martin G. Reese,<sup>96</sup> Knut Reinert,<sup>97</sup> Karin Remington,<sup>98</sup> Robert D. C. Saunders,<sup>99</sup> Frederick Scheeler,<sup>100</sup> Hua Shen,<sup>101</sup> Brian Christopher Shue,<sup>102</sup> Inga Sidién-Kiamos,<sup>103</sup> Michael Simpson,<sup>104</sup> Marian P. Skopaki,<sup>105</sup> Tom Smith,<sup>106</sup> Eugene Spiro,<sup>107</sup> Allan C. Spradling,<sup>108</sup> Mark Stapleton,<sup>109</sup> Renee Strong,<sup>110</sup> Eric Sun,<sup>111</sup> Robert Svirskis,<sup>112</sup> Cyndee Tector,<sup>113</sup> Russell Turner,<sup>114</sup> Eli Venter,<sup>115</sup> Alhui H. Wang,<sup>116</sup> Xin Wang,<sup>117</sup> Zhen-Yuan Wang,<sup>118</sup> David A. Wasserman,<sup>119</sup> George M. Weinstock,<sup>120</sup> Jean Weissenbach,<sup>121</sup> Sherita M. Williams,<sup>122</sup> Trevor Woodage,<sup>123</sup> Kim C. Worley,<sup>124</sup> David Wu,<sup>125</sup> Song Yang,<sup>126</sup> Q. Allison Yao,<sup>127</sup> Jane Ye,<sup>128</sup> Ku-fang Yeh,<sup>129</sup> Jayshree S. Zaveri,<sup>130</sup> Hing Zhan,<sup>131</sup> Guangren Zhang,<sup>132</sup> Qi Zhao,<sup>133</sup> Liangsheng Zheng,<sup>134</sup> Xiangjun Li, Zheng, Fei N. Zhong,<sup>135</sup> Wanyan Zhong,<sup>136</sup> Xiaojin Zhou,<sup>137</sup> Shiqing Zhu,<sup>138</sup> Xiaohong Zhu,<sup>139</sup> Hamilton O. Smith,<sup>140</sup> Richard A. Gibbs,<sup>141</sup> Eugene W. Myers,<sup>142</sup> Gerald M. Rubin,<sup>143</sup> J. Craig Venter<sup>144</sup>

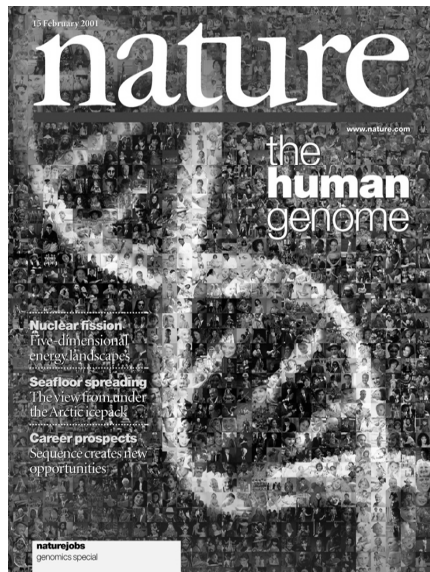
Adams et al. (2000)

# Human Genome Sequencing Centers

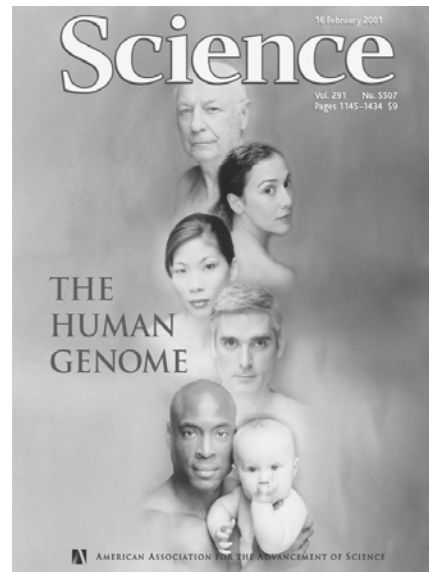




## February, 2001 Draft Sequence



International Human Genome Sequencing Consortium (2001)



Venter et al. (2001)

## April, 2003 Completion



## October, 2004 Publication



**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

21 October 2004

**nature**

www.nature.com/nature

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 291 megabases) contained 2.85 billion nucleotides interrupted by only 2,411 genes. 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 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 human genes including estimates of gene number, birth and death. Notably, the human genome seems 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 high accuracy 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 and (2) the sequencing of the genome. The goal was to obtain a sequence of the genome with sufficient accuracy to allow the study of inherited disease and provide a critical scaffold for genome assembly and (3) the sequencing of repetitive and mobile, single genomes (4) to serve as a scaffold for genetic diagnosis and as a tool for 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 six countries, was formed to carry out this component of the HGP.

In February 2001, the IHGSC and Celera Genomics\* 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 gene clusters, local architecture of genes, regional differences in genome complexity, distribution and history of transposable elements, distribution of polymorphisms and relationships between genetic recombination and physical distance. Moreover, systematic knowledge of the human genome has enabled new tools and approaches that have reshaped academic biomedical research.

Both draft sequences, however, had important shortcomings. The IHGSC sequence, for example, covered ~90% of the euchromatic genome. It was interrupted by ~150,000 gaps and the order and orientation of many sequences within local regions had not been established. The Celera sequence, this turned to the challenge of completing the sequence of the euchromatic genome. Systematically, a finished sequence was defined as having an error rate of, at most, one error per 10<sup>6</sup> bases, and the goal for completion was coverage in finished sequence of at least 95% of the euchromatic genome, with the only gaps being those not amenable to sequencing (see <http://www.genome.gov/1000923>). The goal was challenging because the human genome sequence with such resolution is interrupted with large repetitive duplications, which greatly complicate the determination of genome euchromatic sequence. In fact, near complete sequences have been obtained so far only for two model organisms: the nematode, *Caenorhabditis elegans* and the yeast, *Saccharomyces cerevisiae*. These genomes are, at roughly 10-fold smaller than the human genome and have much simpler structure.

We conclude from the results of a midlayer effort by the IHGSC towards the goal of a complete human sequence. The number of gaps has been reduced 40-fold to only 261, most of which are associated with segmental duplications 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 2.85 billion nucleotides across ~99% of the euchromatic 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 studies made possible by the sequence. We also discuss some accepted broader aspects of the content of the human genome. An initial analysis of the content of the human genome, including a search for previously reported and a set of gaps, is being written describing the structural characteristics (5), including 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 divided in to building (1) a complete, iterative genome that generally encompasses all multiple copies, ranging from single nucleotide to the origin of whole chromosome copies. The fundamental challenge is that genome copies that are not well represented or easily resolved through random shotgun sequencing tend to be highly enriched in problematic sequences. Resolving such regions required the development of special approaches, which evolved substantially over time and varied among centers.

Initially, the finishing process involved two distinct components: (1) producing finished maps, consisting of contigs and/or supercontigs of overlapping large insert clones spanning the euchromatic region of the chromosome arms 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 approaches are described below and 2. Further information about the finishing process and finishing standards can be found in the Supplementary Information (S1) and at <http://www.genome.gov/1000923>.

In total, we generated a shotgun sequence from 84,268 large-insert clones (total length ~1.04 gigabases (Gb)) and finished the sequence from 43,742 of these clones (total length ~1.67 Gb). The clones consisted primarily of bacterial artificial chromosomes

## 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.

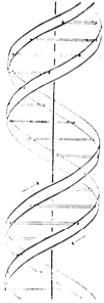
Two men from two separate groups -- Francis Collins of the National Institutes of Health and Craig Venter of Celera Genomics Inc., a pharmaceutical-development company -- worked independently to discover the sequence of the human genome and identify the genes that it contains. This



**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

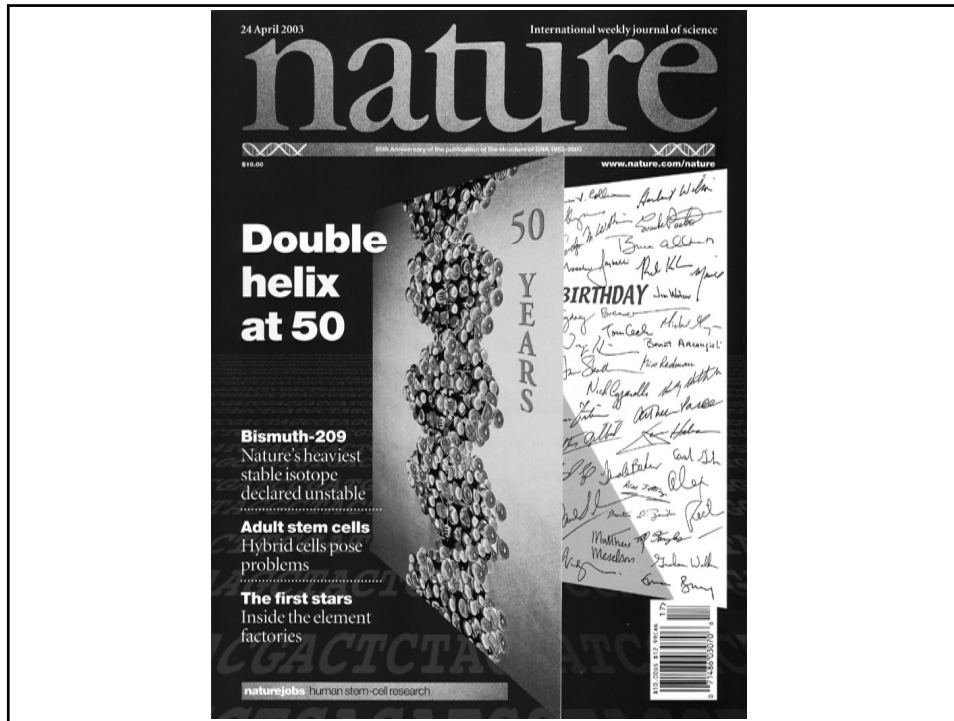
Medical Research Council Unit for the Study of the Molecular Structure of Biological Systems,  
Cavendish Laboratory, Cambridge.  
April 2.



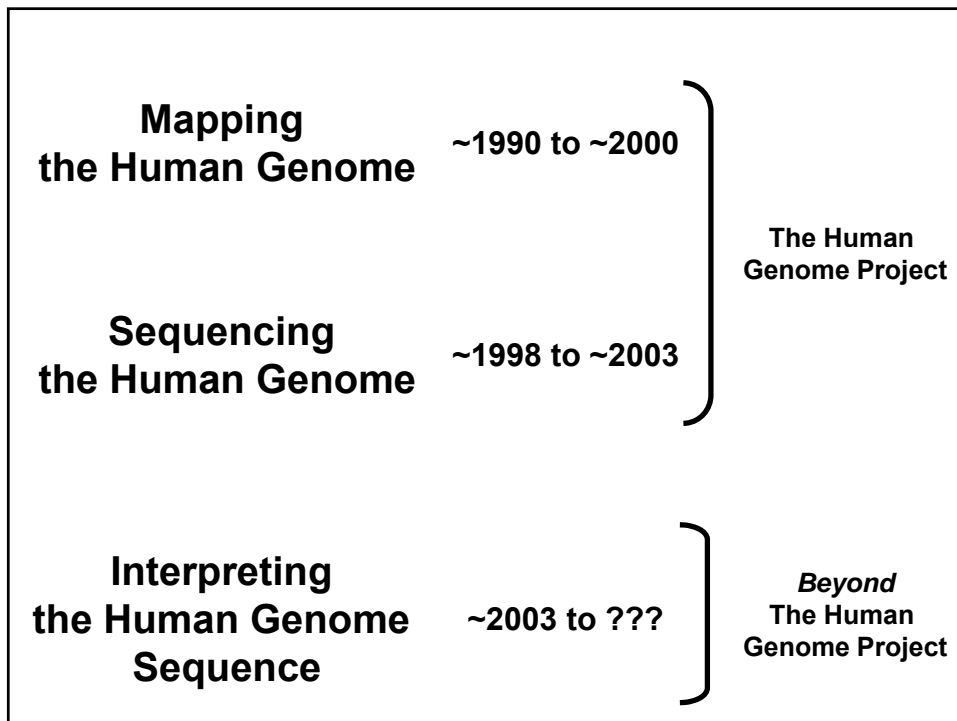
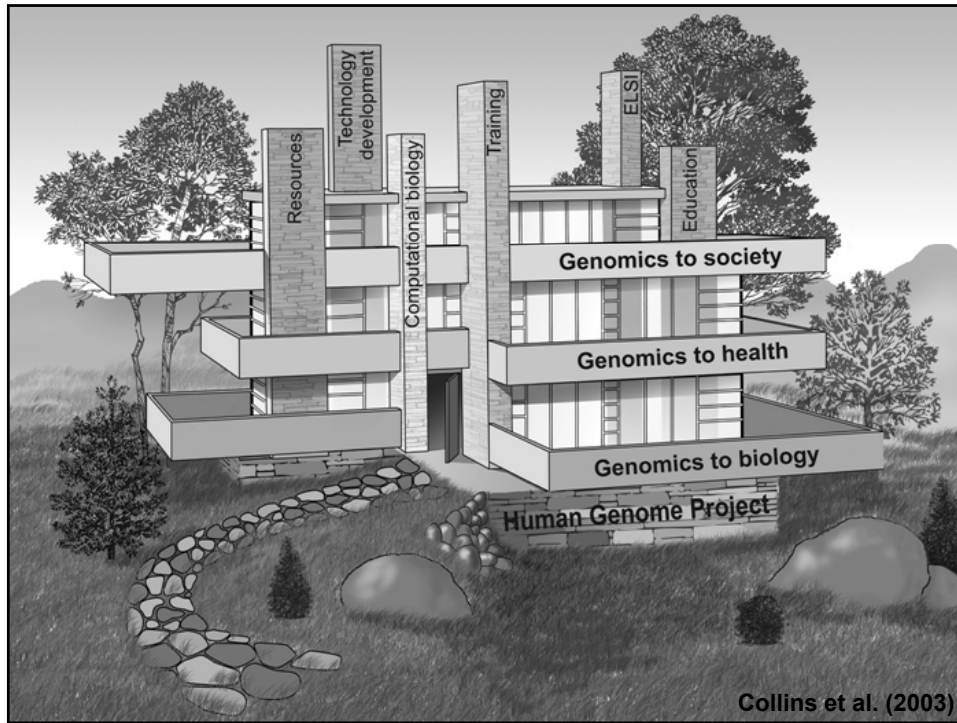
**DOUBLE  
HELIX  
TO  
HUMAN  
SEQUENCE**

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

**What's Next?**



Collins et al. (2003)





## 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**



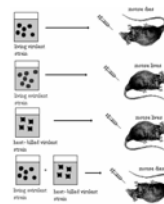
**Mendel**

**1865**



**Miescher**

**1871**



**Avery**

**1944**



**Watson & Crick**

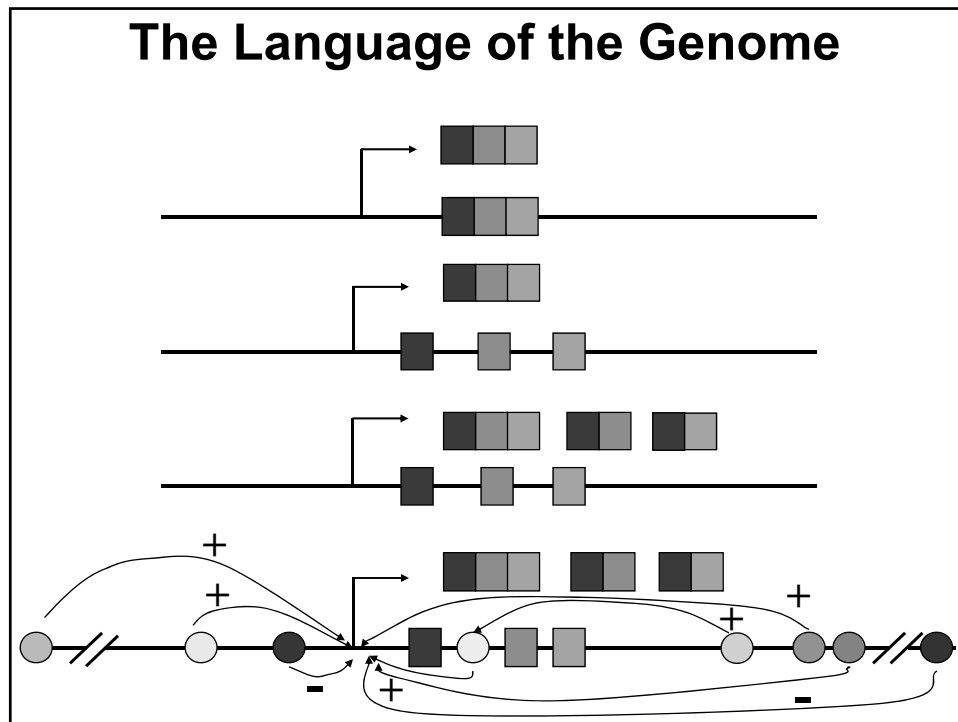
**1953**

## Comparing Genomes is Like Cryptography

CKQEBHEREYTWASUISZMEISDFOGETHEBLPBGODFQSTLKSTUFFRTAC  
 DLUCEHEREZBRTTOISAWNDCDARJJPTHERROFGODERGHCLSTUFFBRHA

## Functional Elements: Coding vs. Non-Coding

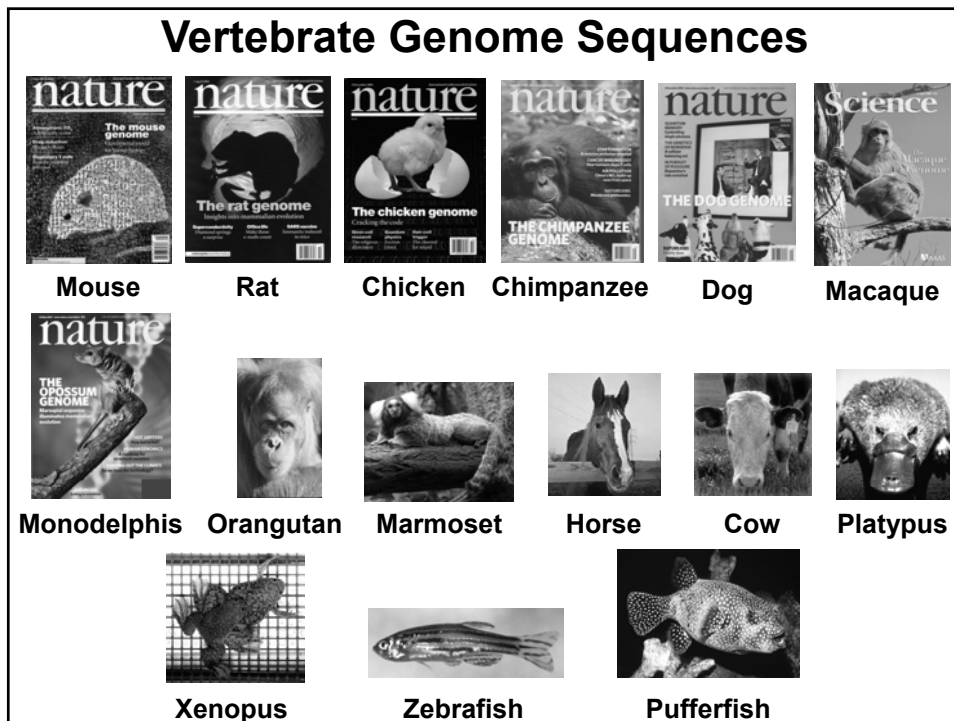
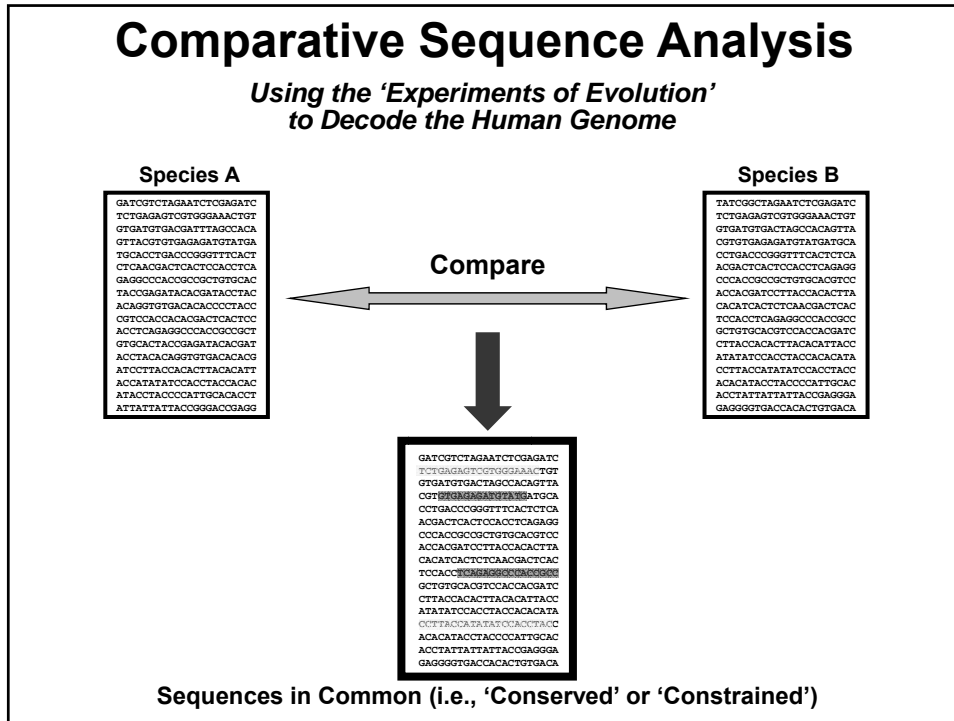
- **Coding Sequences (i.e., Genes)**
  - Relatively EASY to Identify
  - Mostly Know What to Look For
  - Complementary Data Sets Available (ESTs, cDNAs)
  - Ever-Improving Computational Gene Predictions
- **Non-Coding Functional Sequences**
  - HARD to Identify
  - Very Little Known About What to Look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



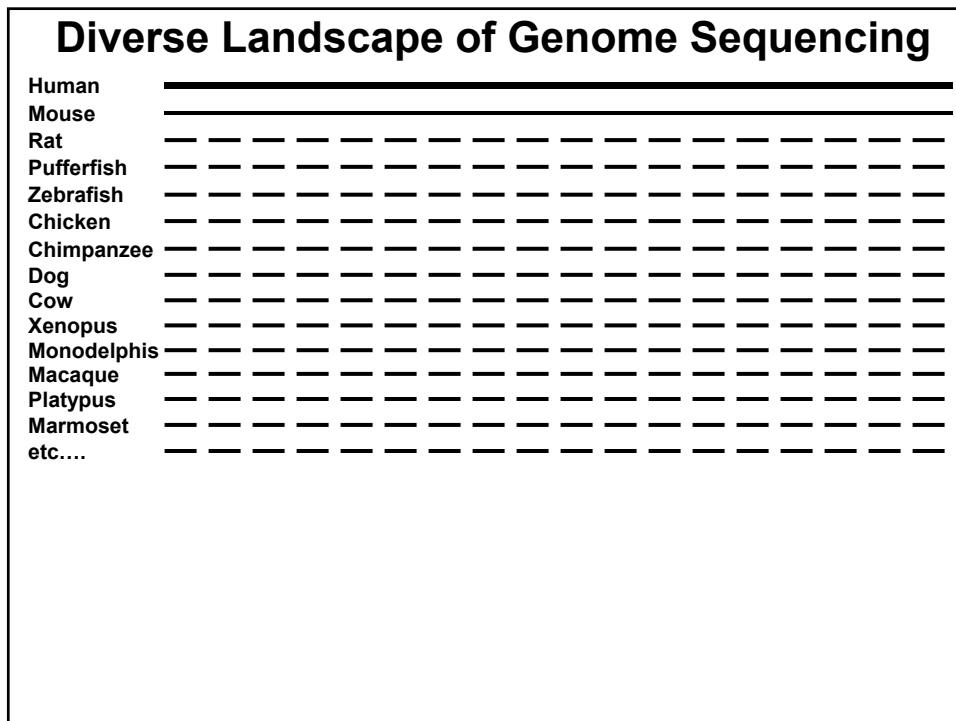
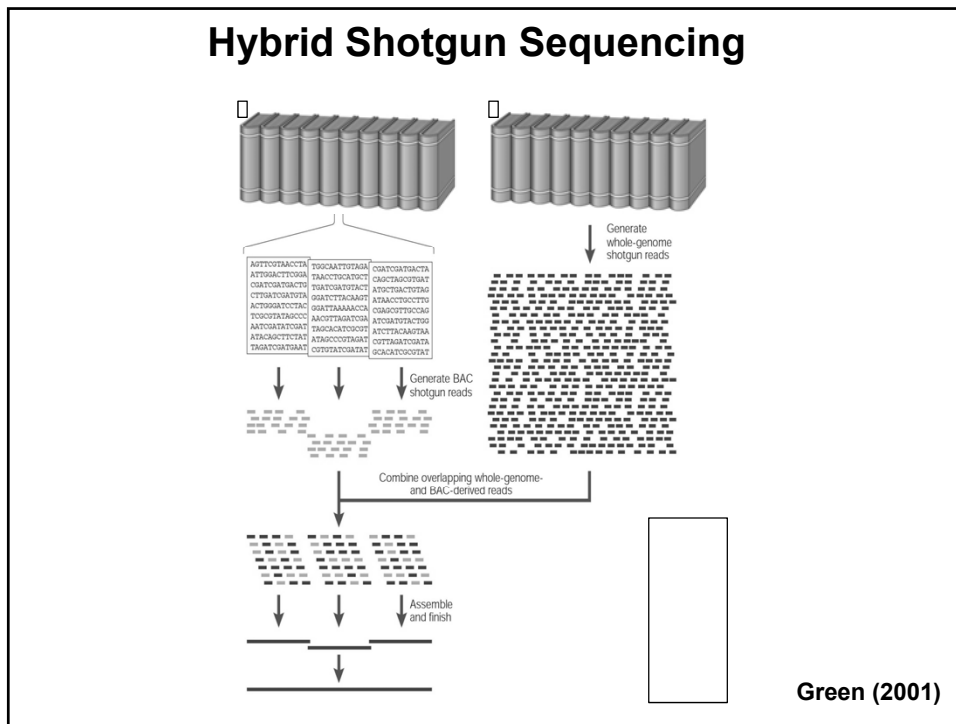
## Functional Elements: Coding vs. Non-Coding

- **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**





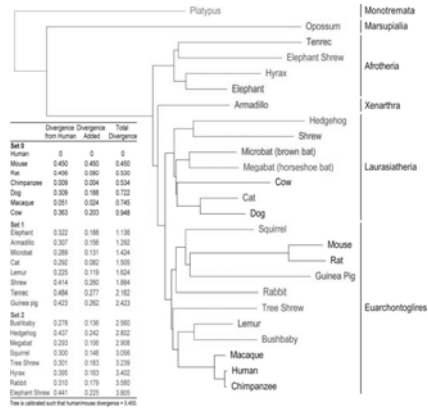


## Low-Redundancy, Whole-Genome Shotgun Sequencing

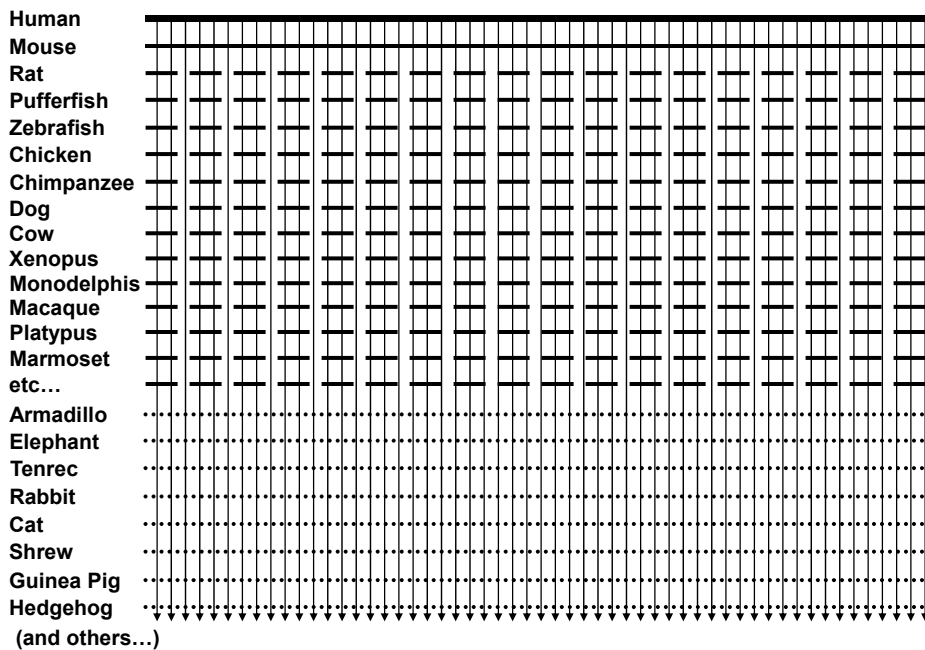
An initial strategy for the systematic identification of functional elements in the human genome by low-redundancy comparative sequencing

Filiott H. Margulies<sup>\*†</sup>, Lada Vinson<sup>†‡</sup>, NISC Comparative Sequencing Program<sup>\*§¶</sup>, Wehh Miller<sup>‡</sup>, David R. Jaffe<sup>‡</sup>, Kerstin Lindblad-Toh<sup>‡</sup>, Jean Chang<sup>‡</sup>, Eric D. Green<sup>\*§</sup>, Eric S. Lander<sup>‡</sup>, James C. Mullikin<sup>\*§\*\*</sup>, and Michele Clamp<sup>\*\*</sup>

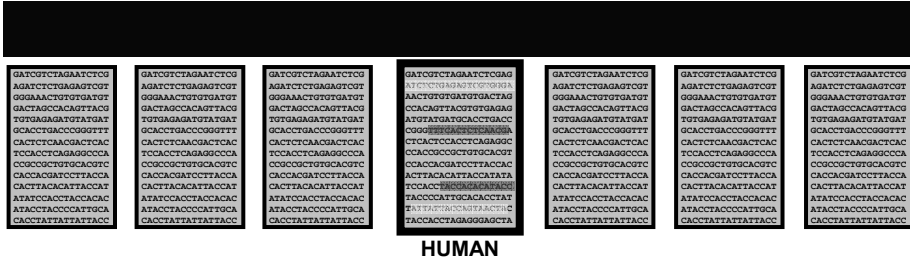
Margulies et al. (2005)



## Diverse Landscape of Genome Sequencing



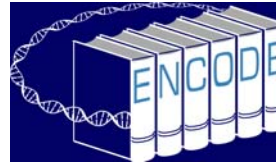
## Multi-Species Sequence Comparisons



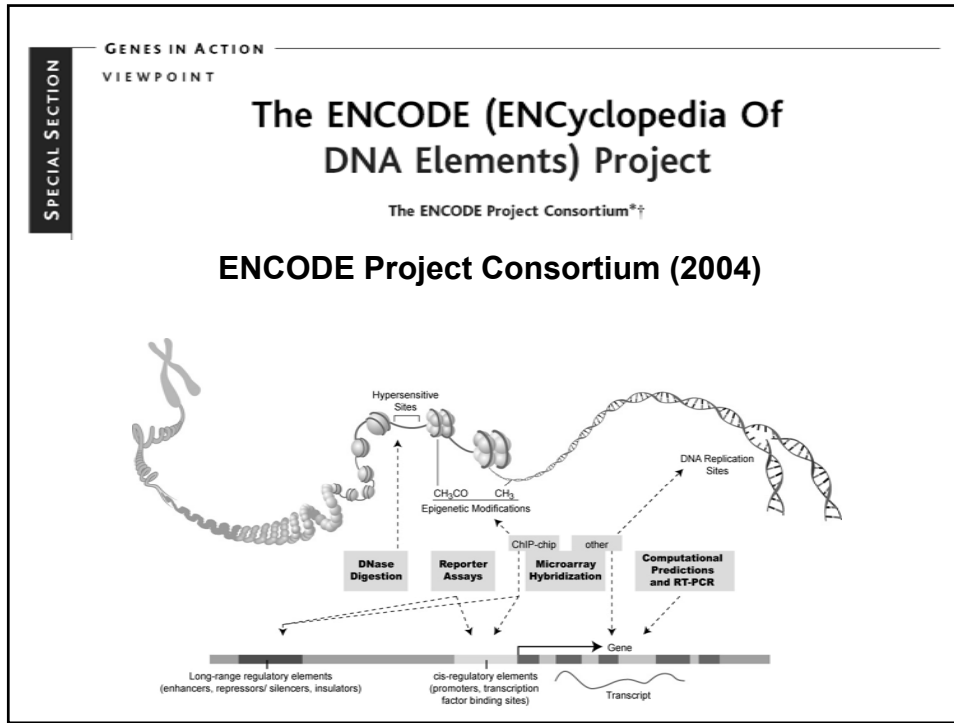
### Multi-Species Conserved Sequences (MCSs)

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

## ENCODE Project



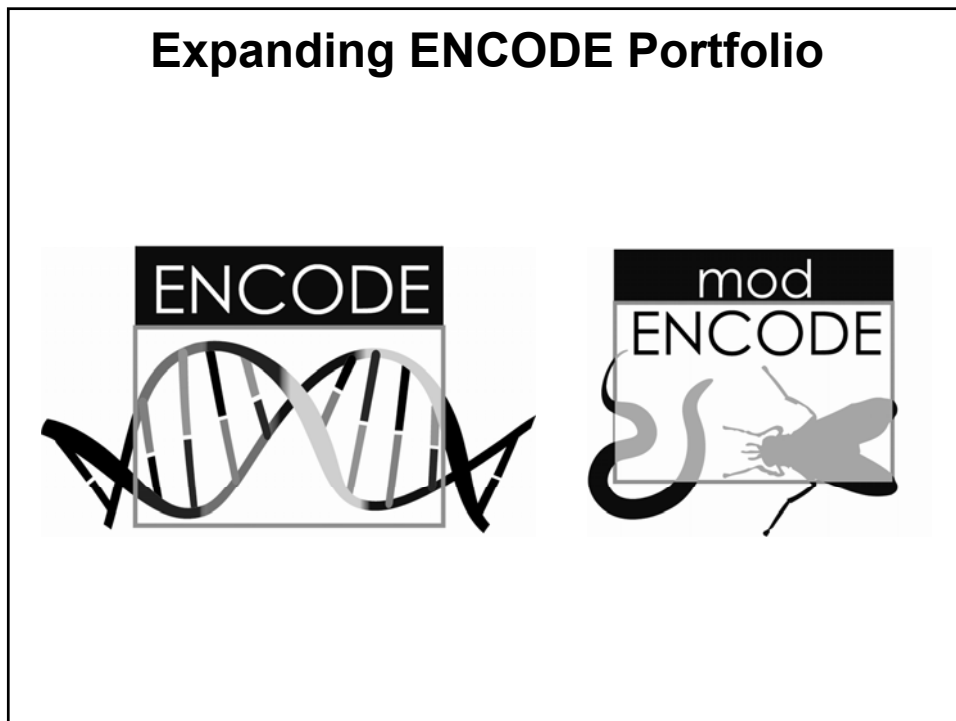
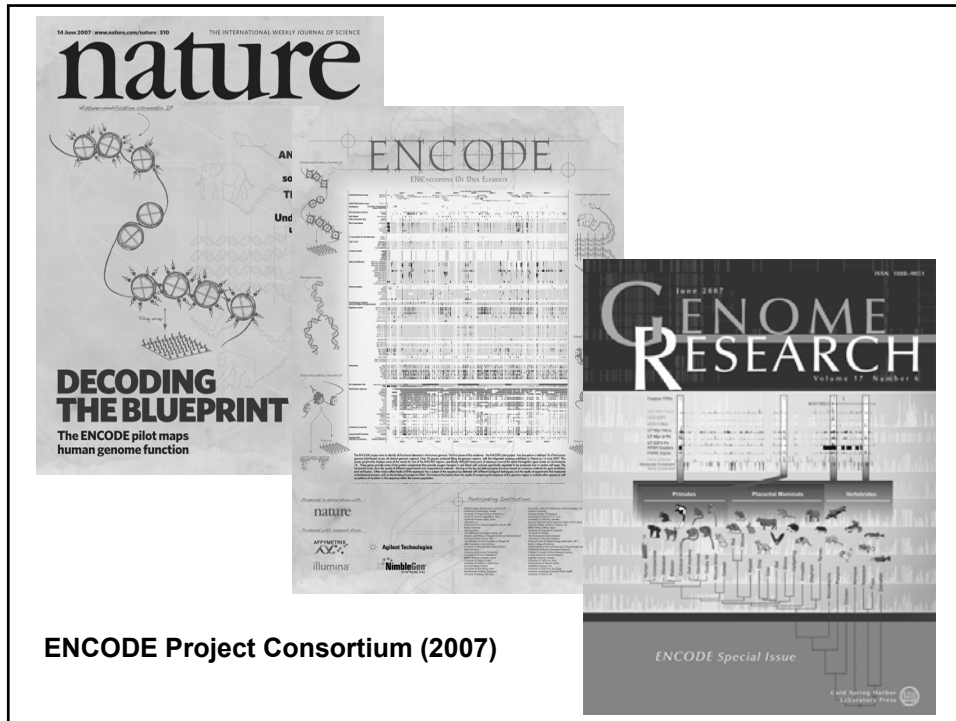
- ENCODE: ENCyclopedia Of DNA Elements
- Goal: Compile a *Comprehensive Encyclopedia of All Functional Elements in the Human Genome*
- Initial Pilot Project: 1% of Human Genome
- Apply Multiple, Diverse Approaches to Study and Analyze that 1% in a Consortium Fashion

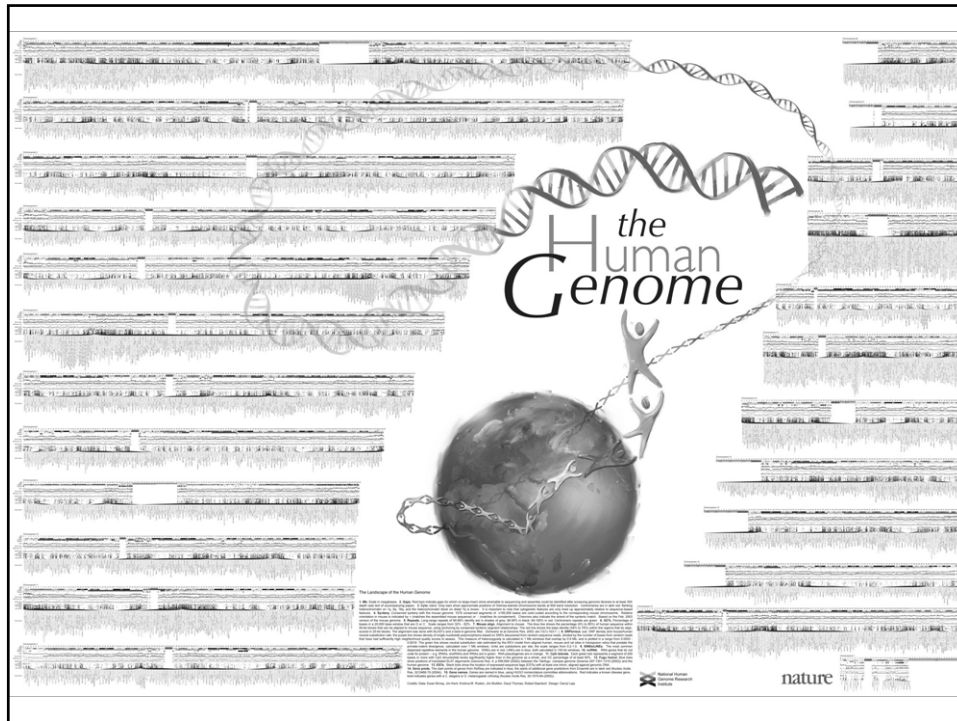


## ENCODE Project: Web Sites

**genome.gov/ENCODE**

**genome.ucsc.edu/ENCODE**

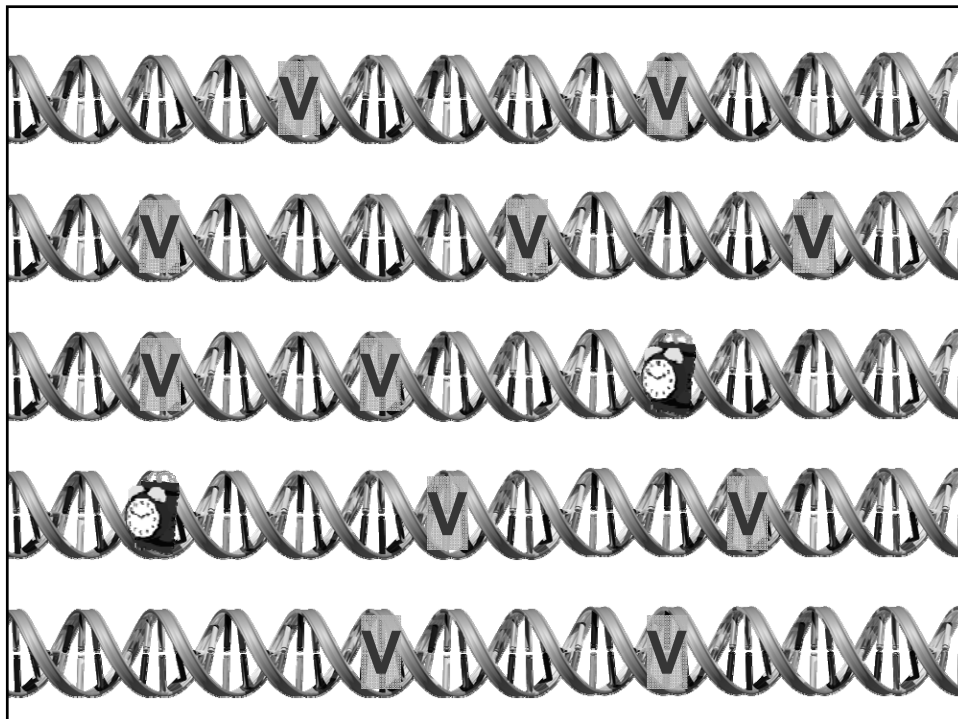


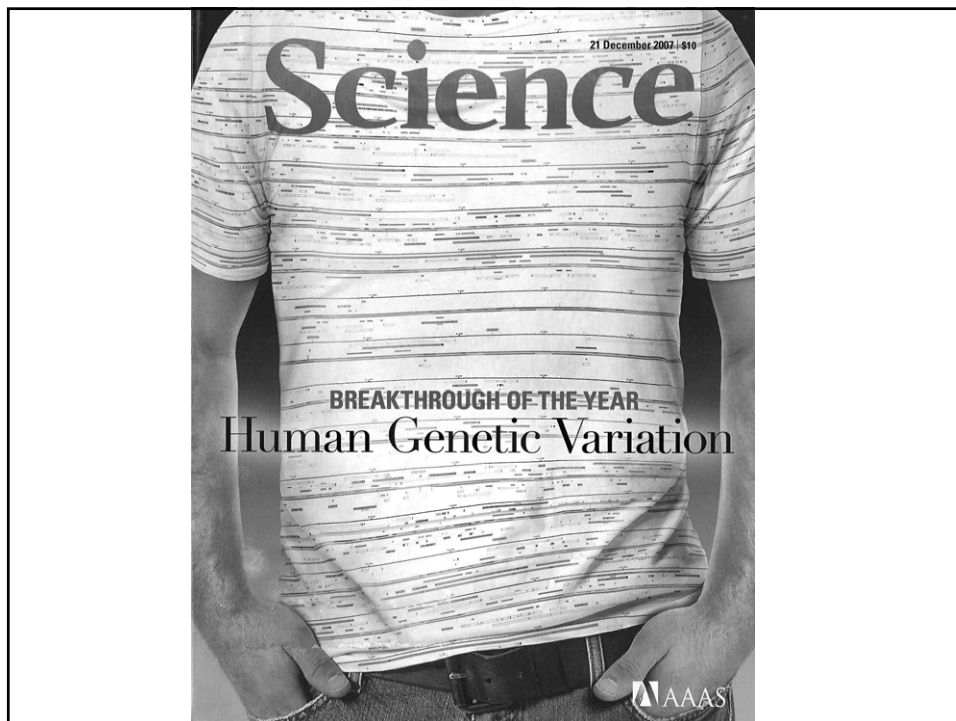


**All humans are ~99.9% identical at the DNA sequence level, and yet...**



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









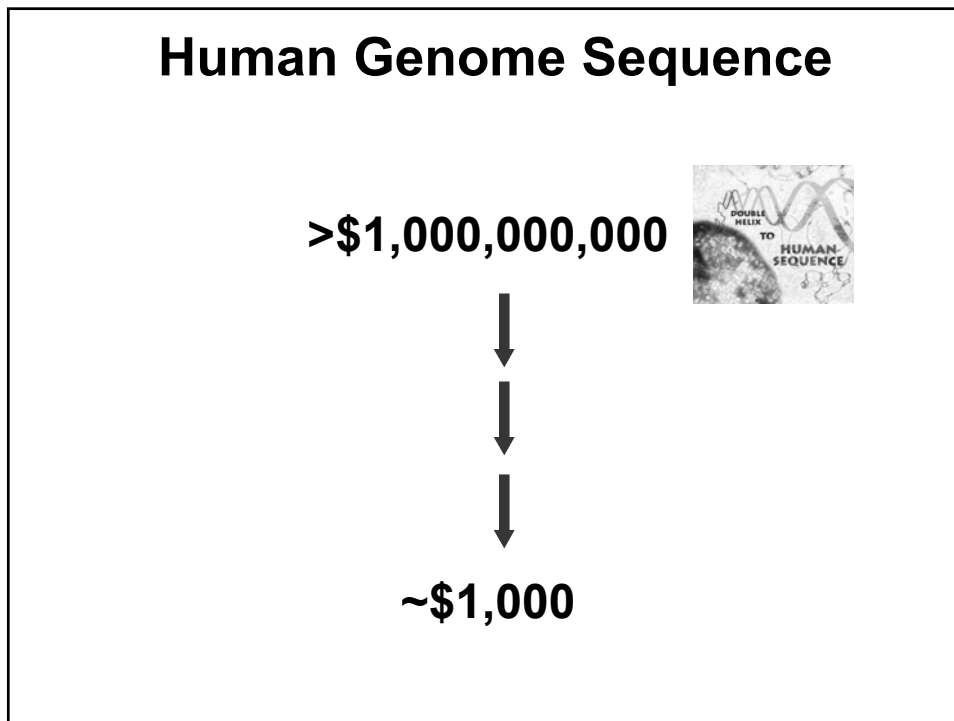
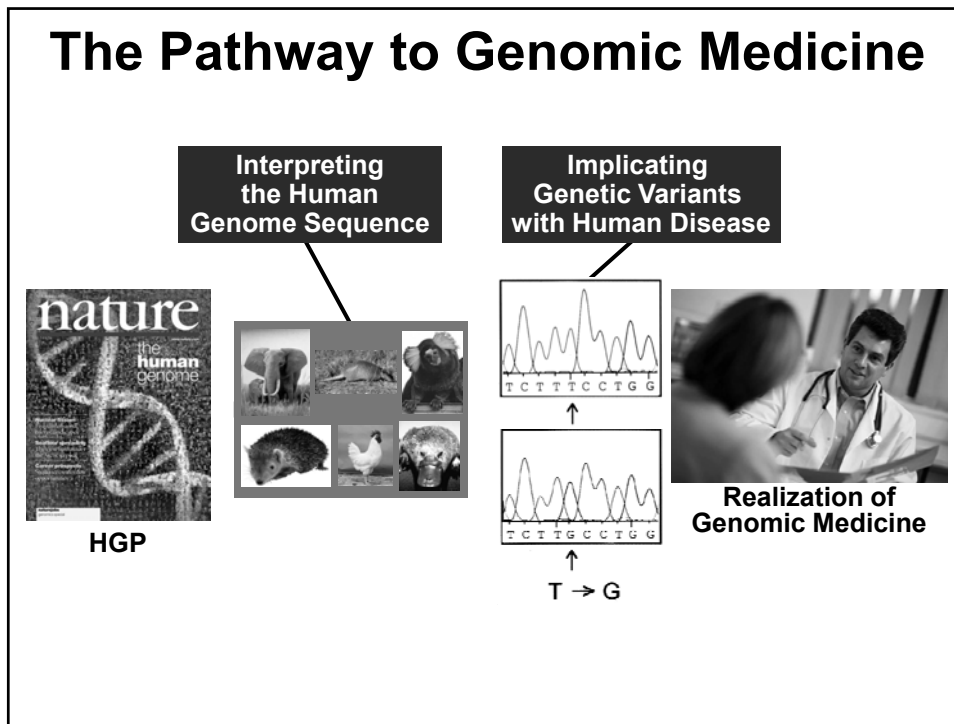
## The Pathway to Genomic Medicine

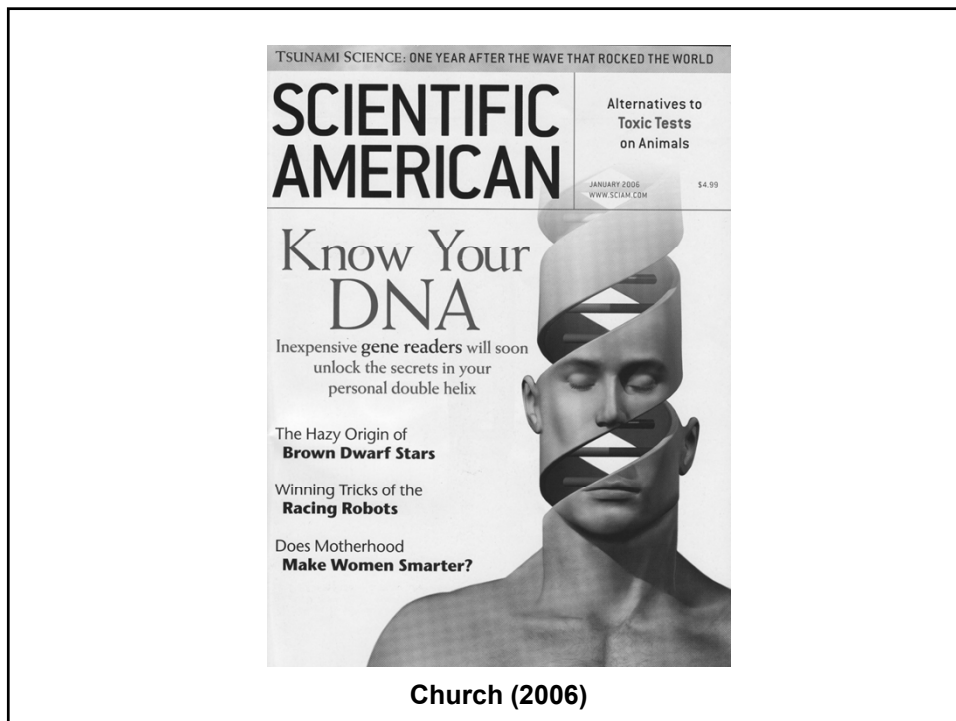
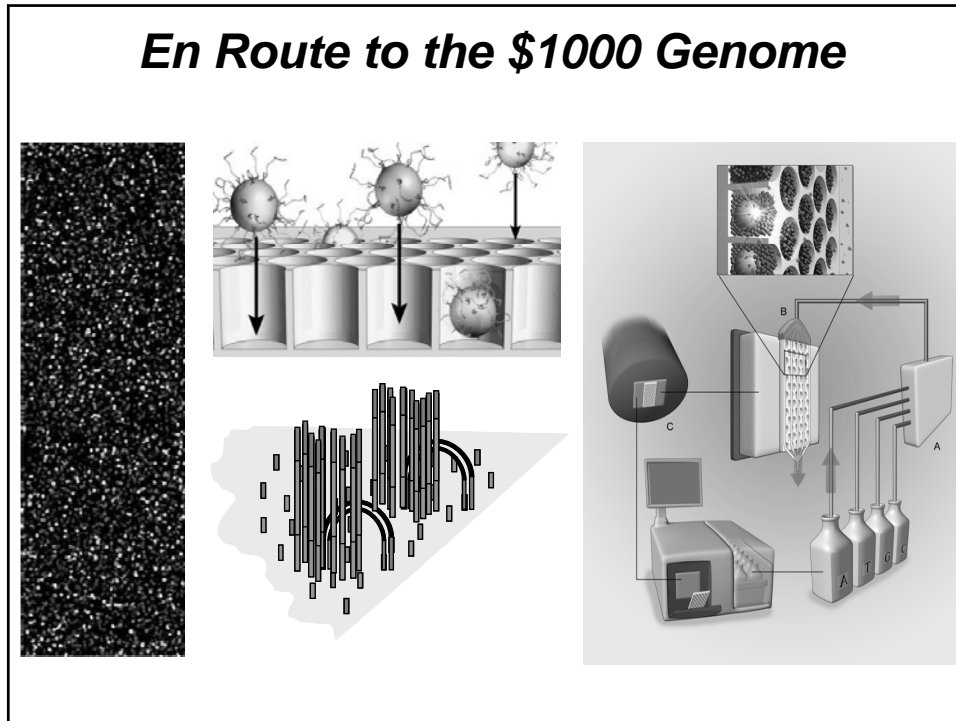


HGP



Realization of Genomic Medicine





**Realities of New DNA Sequencing Technologies...**

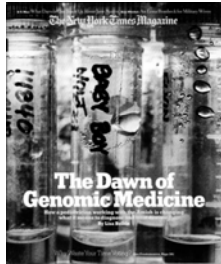


**Changing Infrastructure Requirements**

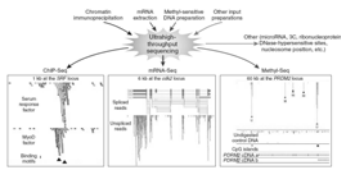
**DNA Sequence  
Production**

**Bioinformatic  
Analysis**

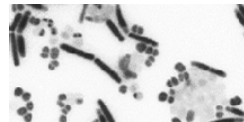
### Expanding Universe of Sequence-Based Explorations



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.