




National Human Genome Research Institute Division of Intramural Research



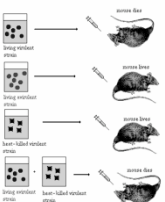
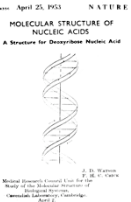


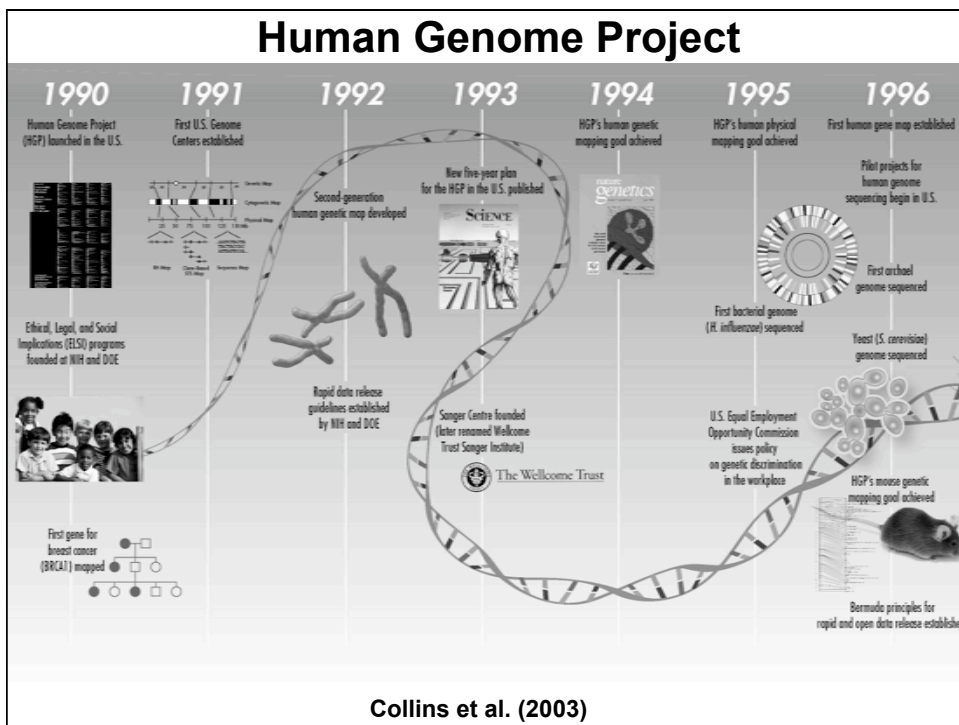
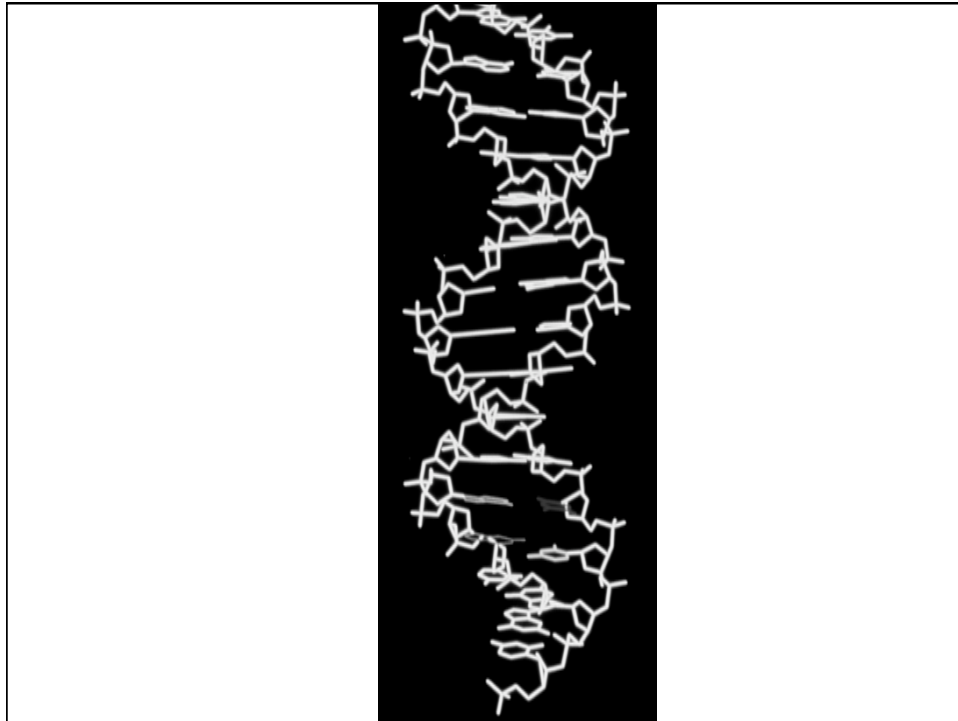
**The Genomic Landscape:
circa 2010**

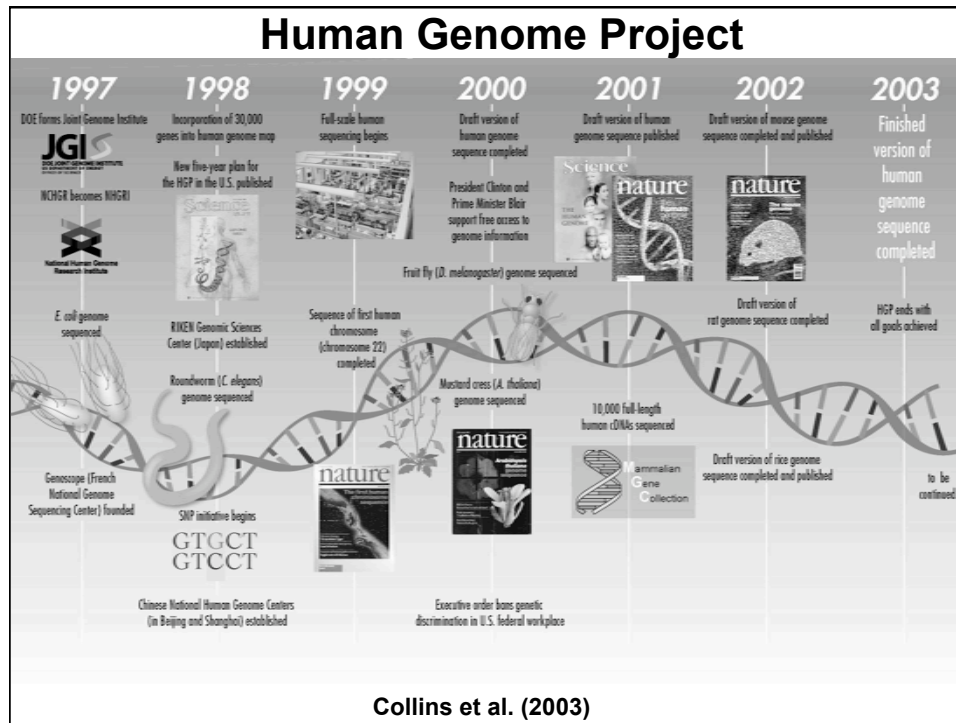
**Eric Green, M.D., Ph.D.
Director, NHGRI**



Foundational Milestones in Genetics & Genomics

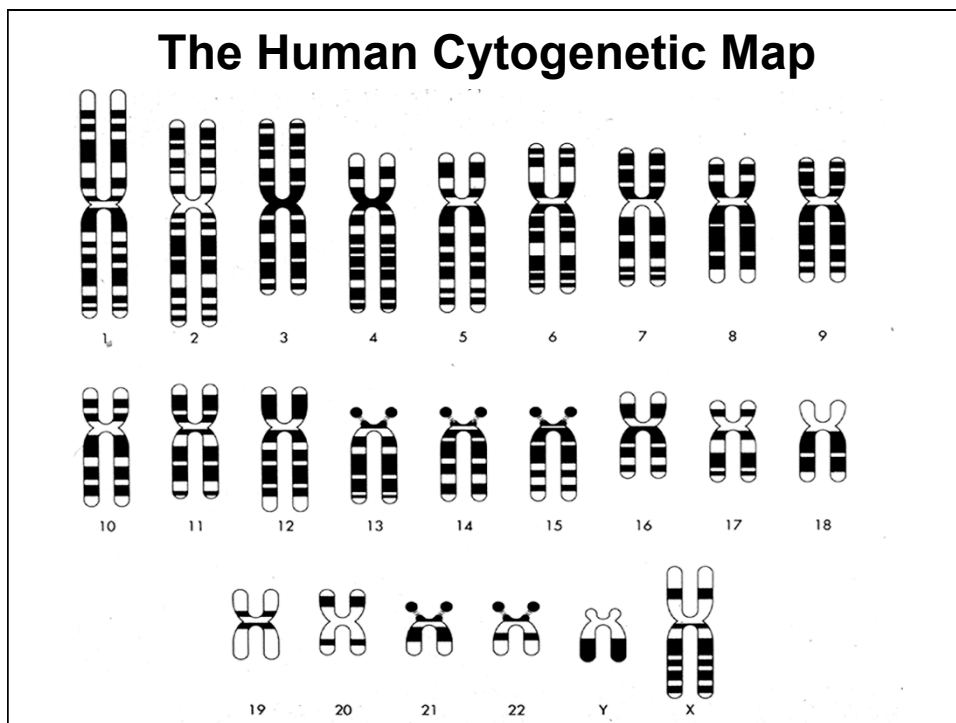
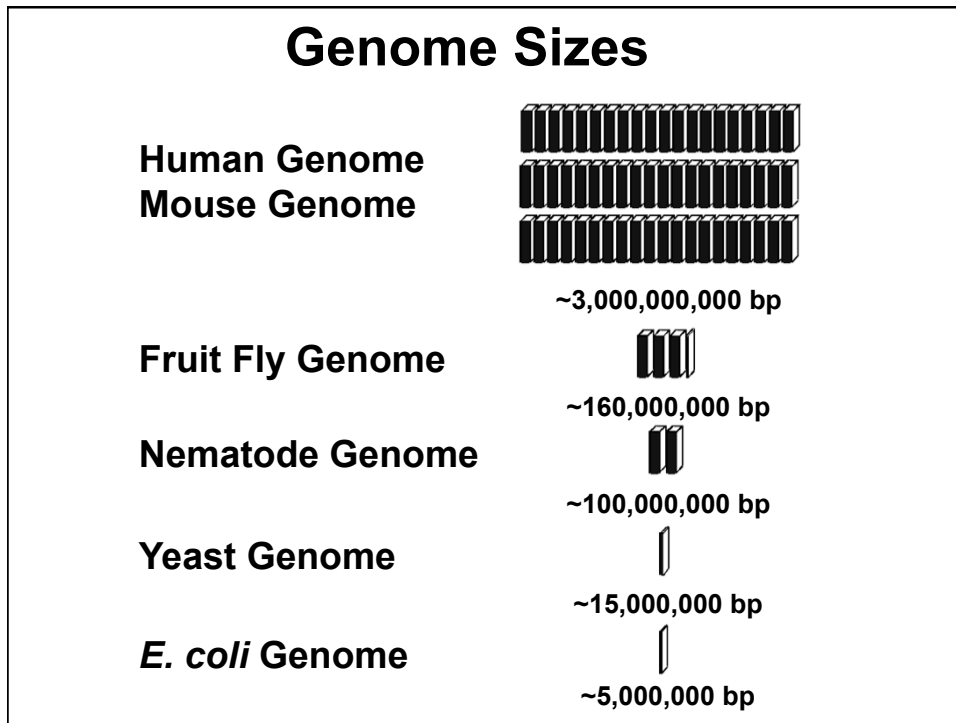
			
Mendel	Miescher	Avery	Watson & Crick
1865	1871	1944	1953

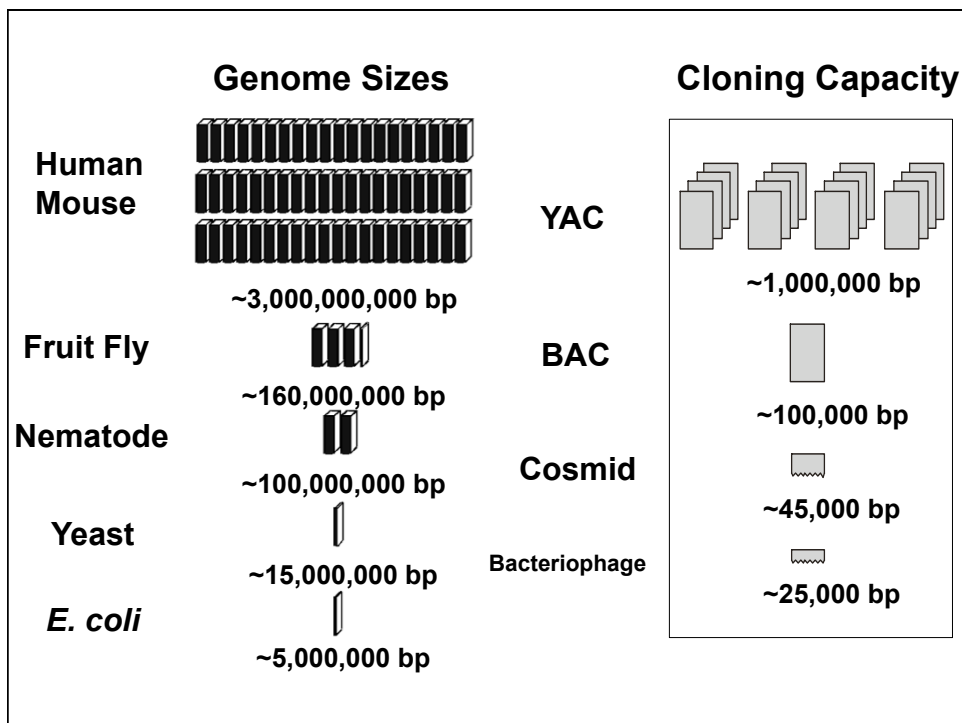
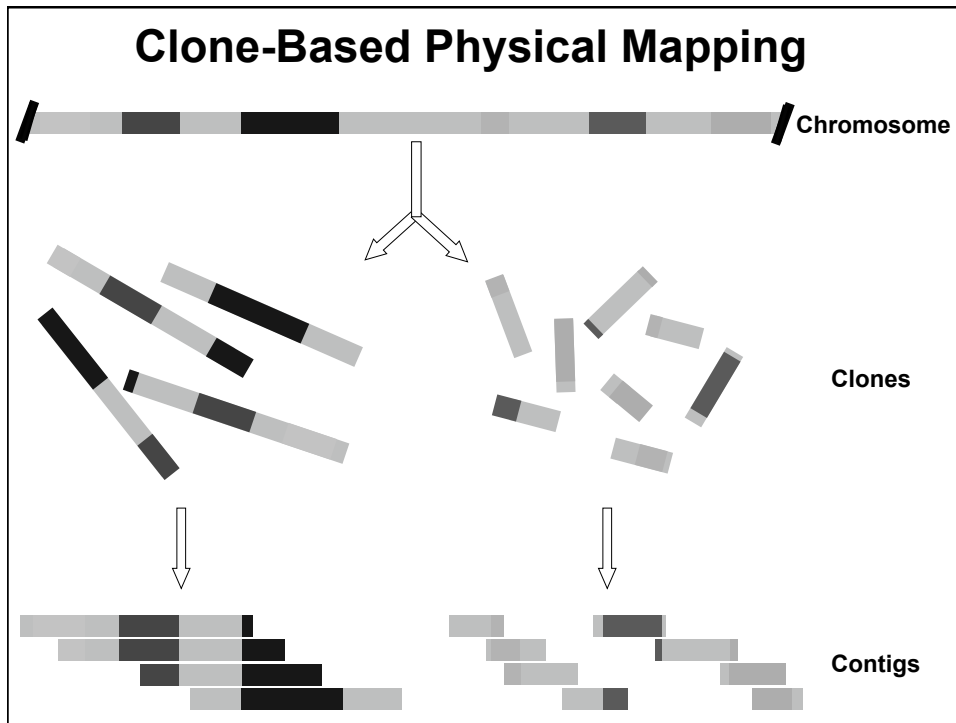


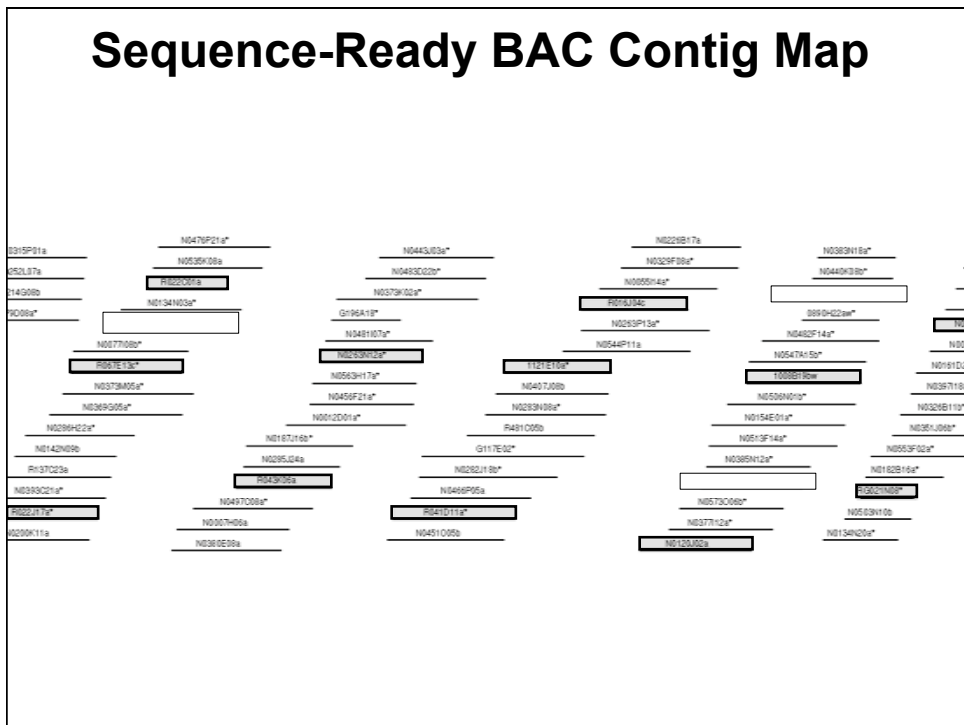
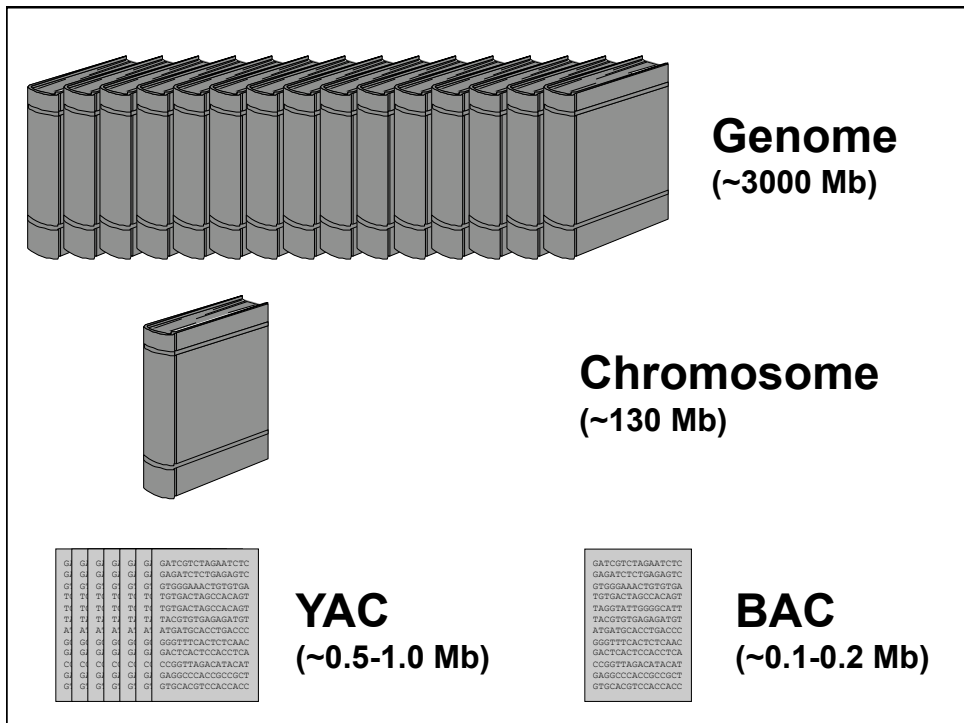


Outline

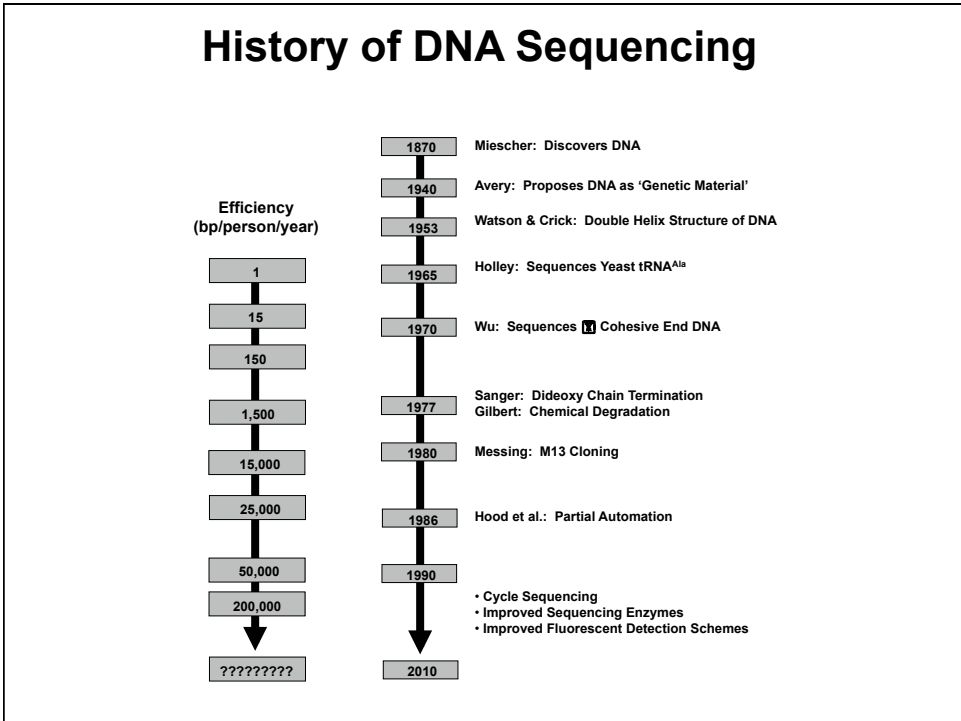
- I. Fundamentals of Genome Mapping & Sequencing
- II. Mapping & Sequencing in the Human Genome Project
- III. Comparative Sequencing
- IV. New Frontiers in Genomics



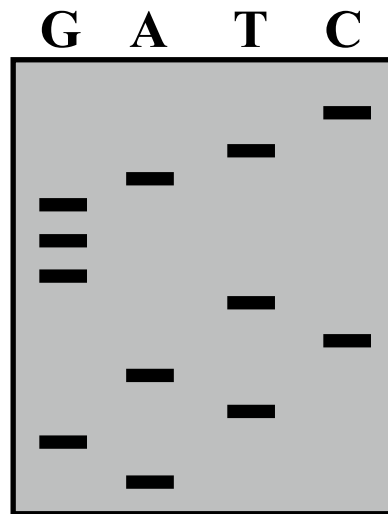




DNA Sequencing

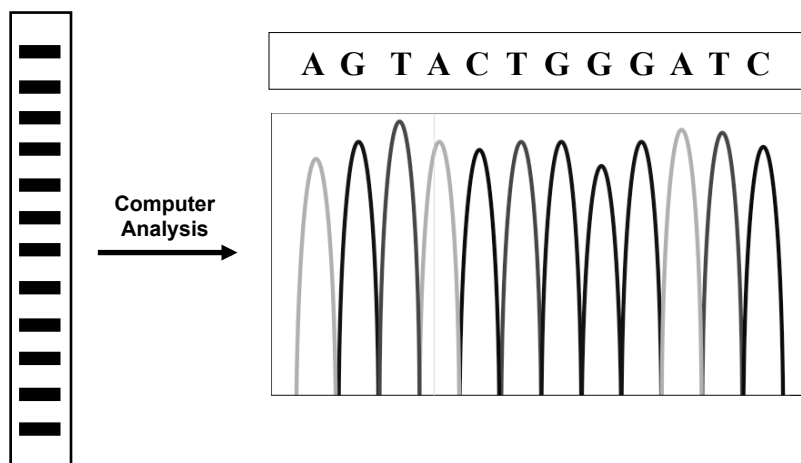


DNA Tagged with Radioactivity

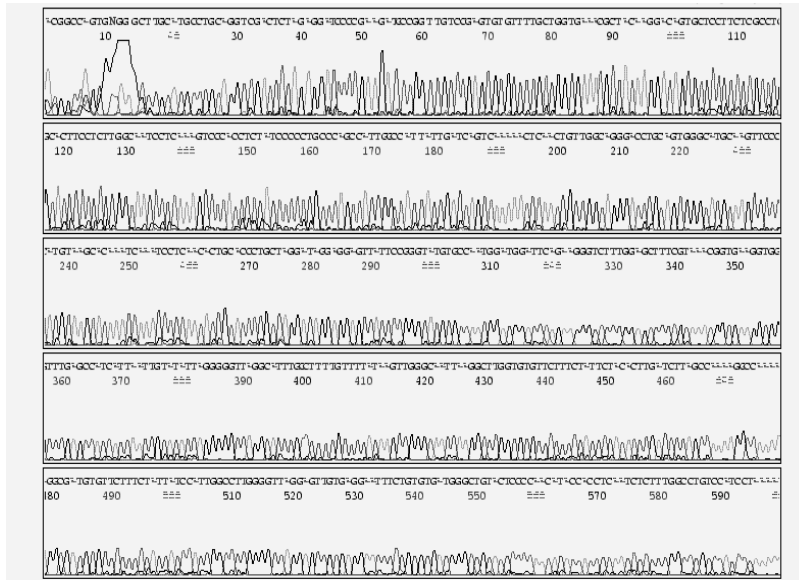


G: G Reaction
A: A Reaction
T: T Reaction
C: C Reaction

Analyzing Fluorescent DNA Sequencing Data

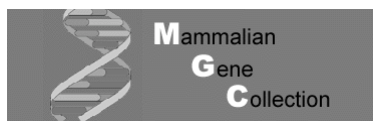


Fluorescent DNA Sequencing Results

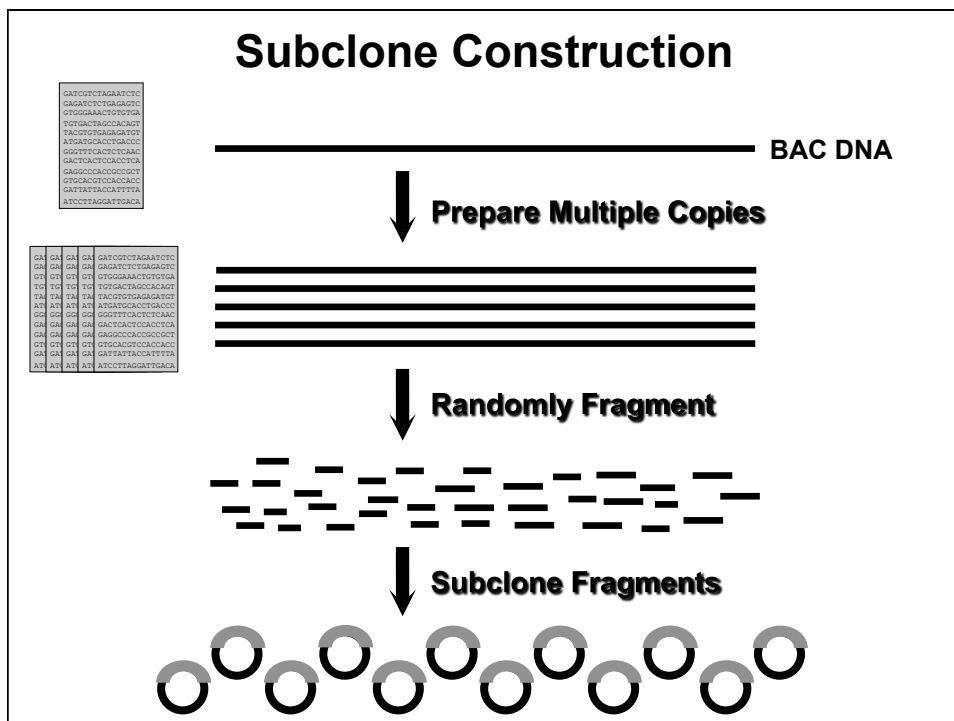
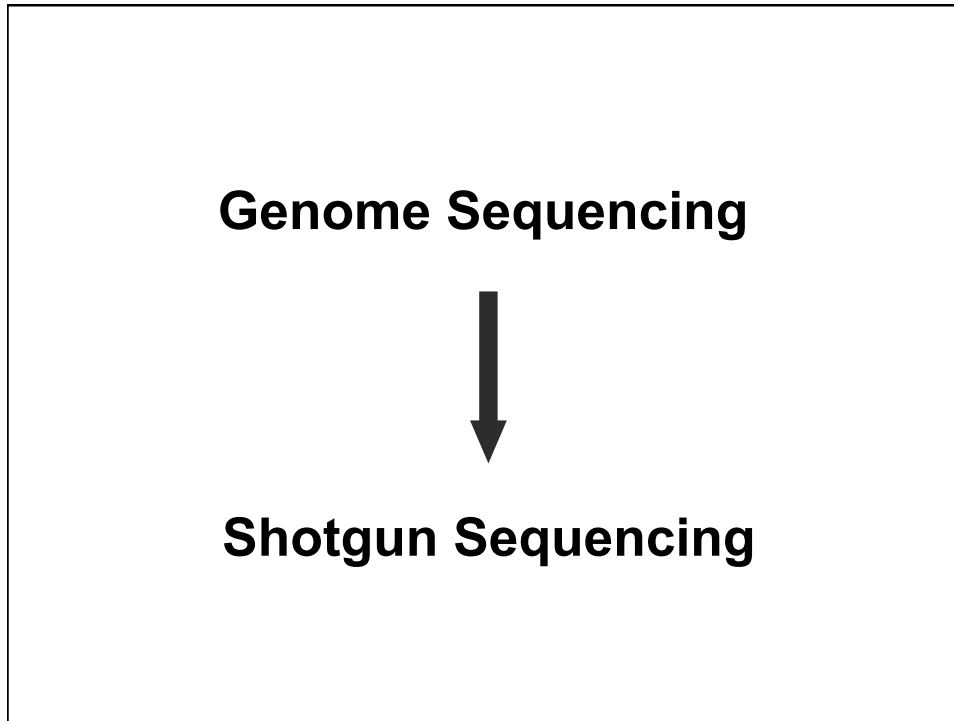


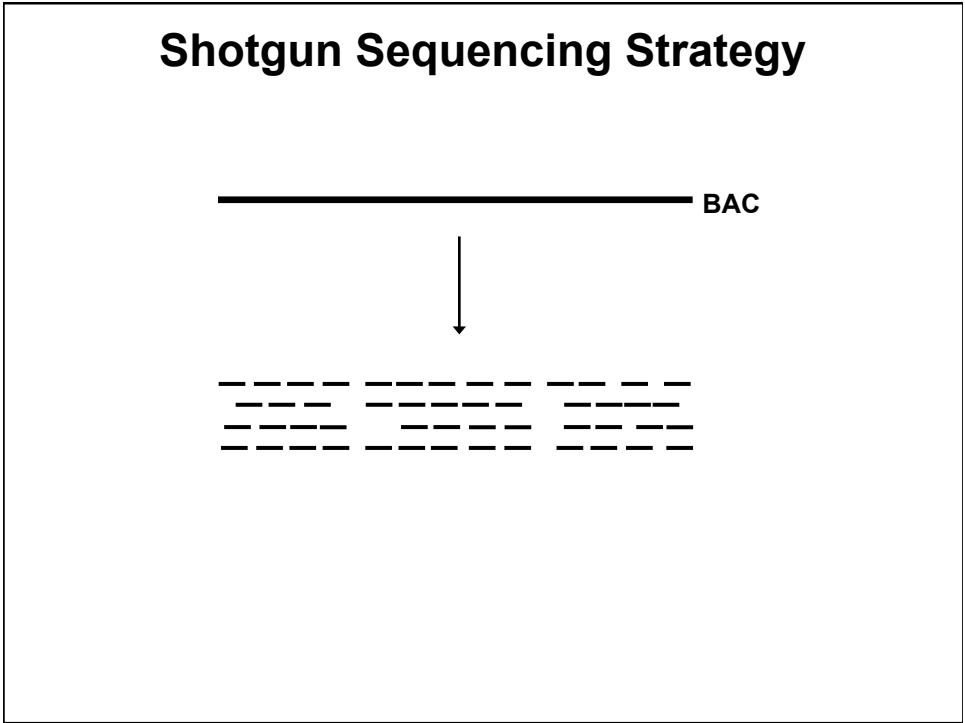
Analysis of Gene Expression

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



mgc.nci.nih.gov





Poisson Calculations

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

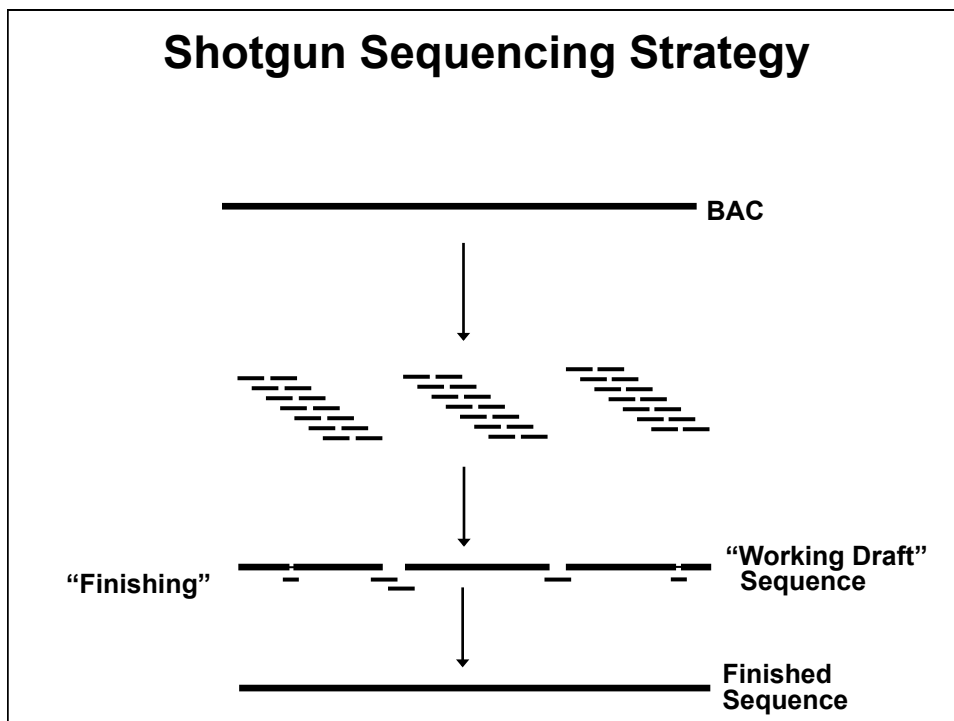
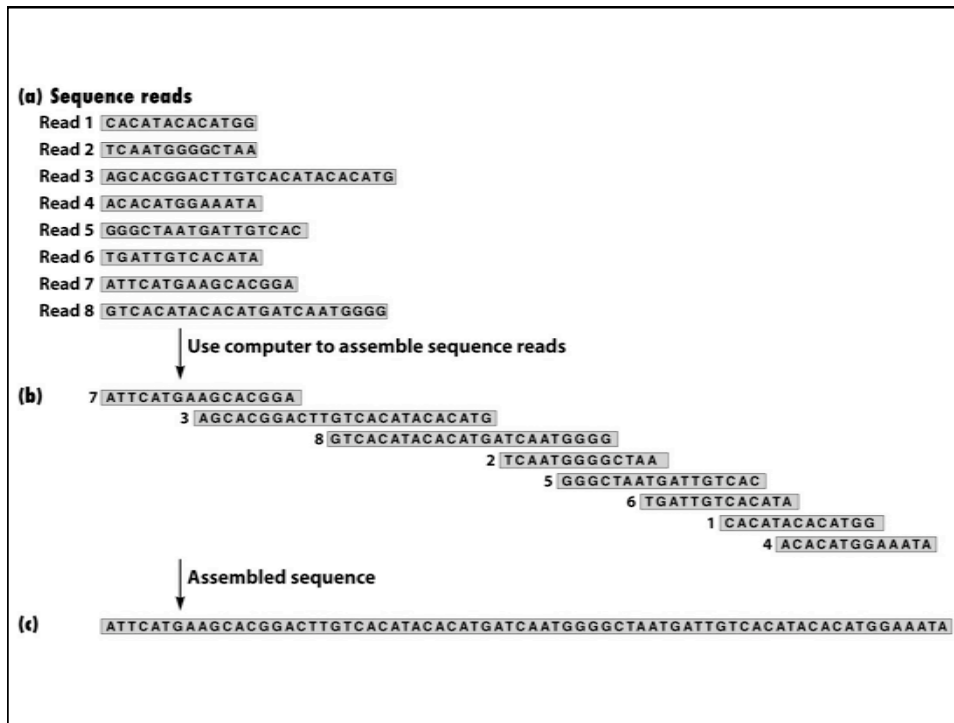
The probability a base is not sequenced is given by:

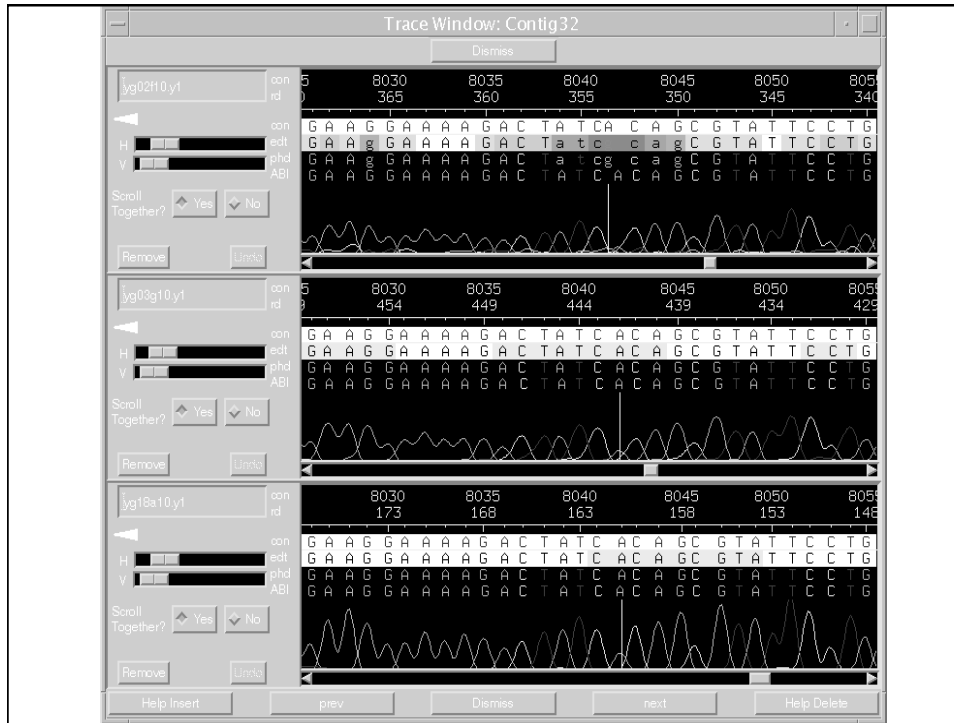
$$P_0 = e^{-c}$$

Where:

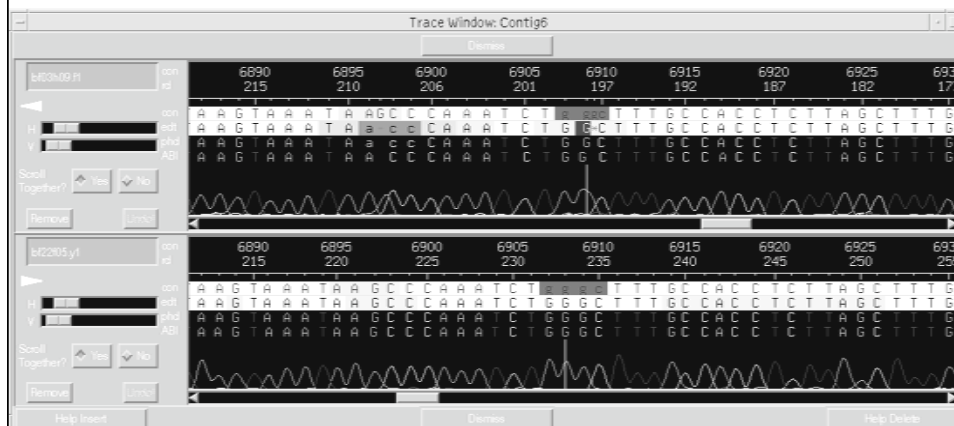
- c = fold sequence coverage ($c = LN/G$),
- LN = # bases sequenced, i.e. L = average sequencing read length and N = # reads
- G = target sequence length
- $e = 2.718$ ($e = 2.718281828459$)

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%





Sequence Finishing: Resolving Ambiguities



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

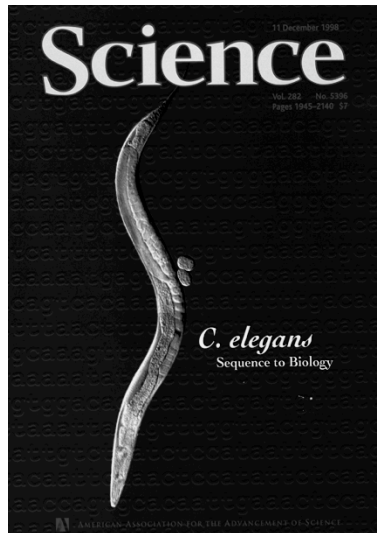
Historically Significant Genome Sequencing Projects

First Eukaryotic Genome Sequence



Goffeau et al. (1997)

First Animal Genome Sequence

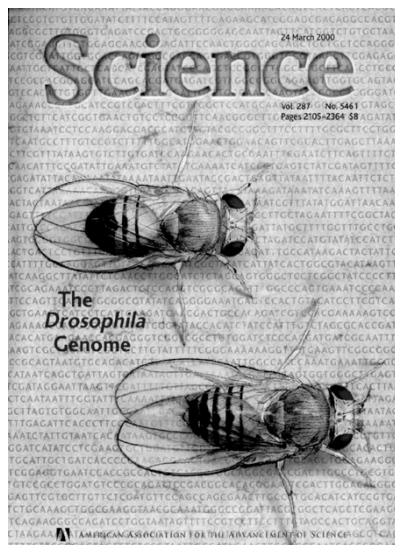


Genome Sequence of the Nematode *C. elegans*: A Platform for Investigating Biology

The *C. elegans* Sequencing Consortium*

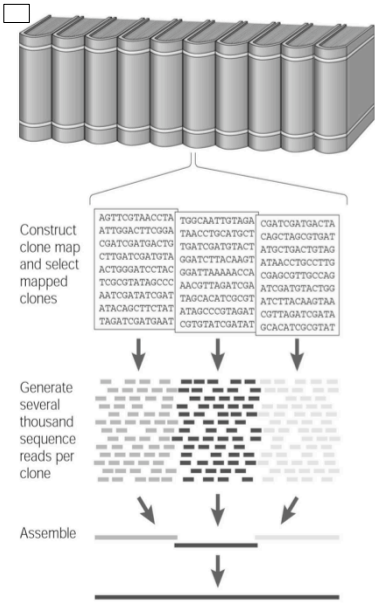
C. elegans Sequencing Consortium (1998)

Second Animal Genome Sequence

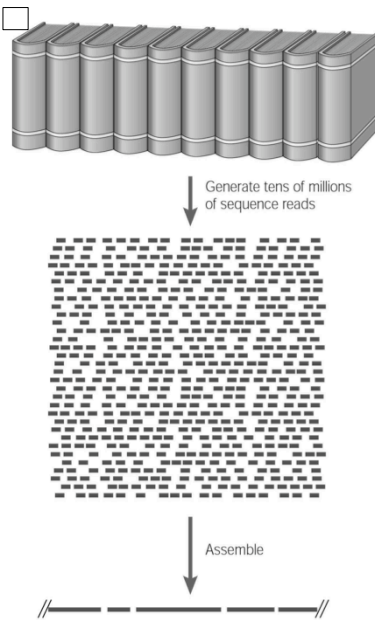


Adams et al. (2000)

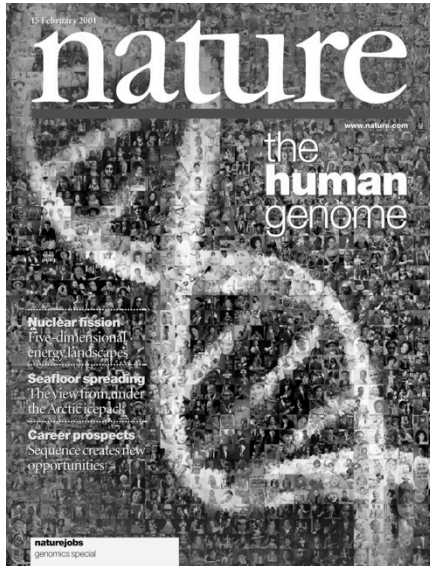
Clone-Based Shotgun Sequencing



Whole-Genome Shotgun Sequencing



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 correct this draft into a genome sequence with high accuracy and nearly complete coverage. Here, we report the result of this finishing process. The current genome sequence (total 2.85 billion nucleotides) comprises 99.99% of the euchromatic genome (3.1 billion nucleotides) and is accurate to an error rate of ~1 error per 100,000 bases. Many of the remaining euchromatic gaps are associated with segmental duplications and are being closed with new methods. The non-coding 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. Notably, the human genome encodes only 28,000–32,000 protein-coding genes. The genome sequence report 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 vast 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 allow the study of inherited disease and provide a scaffold for genome assembly; and (2) the sequencing of sequences with smaller, simpler genomes^{2–4} to serve as a method to 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 (IHSC), an open collaboration involving twenty centres in six countries, was formed to carry out a component of the HGP⁵.

In February 2001, the IHSC and Celis Genomics⁶ each reported draft sequences providing a first overall view of the human genome. These sequences derived from a systematic study of the human genome itself, including identification of genes, combinatorial architecture of proteins, regional differences in genome composition, distribution and history of transposable elements, distribution of polymorphisms and relationships between genes, recombination and physical distance. Moreover, systematic knowledge of the human genome has enabled new tools and approaches that have markedly accelerated biomedical research.

Both draft sequences, however, had important shortcomings. The IHSC sequence, for example, omitted ~10% of the euchromatic genome, was interrupted by ~150,000 gaps and the order and orientation of many segments within local regions had not been established. The IHSC also noted the challenge of completing the sequence of the euchromatic genome. Openly available a finished sequence was defined as having an error rate of at most one error per 10⁷ bases, and the goal for completion was coverage in finished sequence of at least 99% of the euchromatic genome, with the only gaps being those refractory to all available techniques⁷ (see <http://www.genome.gov/090922>). The goal was challenging because the human genome is riddled with such features as dispersed repeats and large segmental duplications, which greatly complicate the determination of genome structure and sequence. In fact, some complete sequences have been obtained so far only for three mitochondrial organisms: the nematode, 'uncultured vertebrate' and the butterfly⁸. These genomes are all roughly 30-fold smaller than the human genome and have much simpler structure.

We describe here the results of a multiyear effort by the IHSC towards the goal of a complete human sequence. The number of gaps has been reduced 400-fold to only 341, 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, a constant 2.85 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 analyses 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 reported⁹ and a series of papers is being written describing the individual chromosomes^{10–12}, including annotation of genes and other features.

Current genome sequence

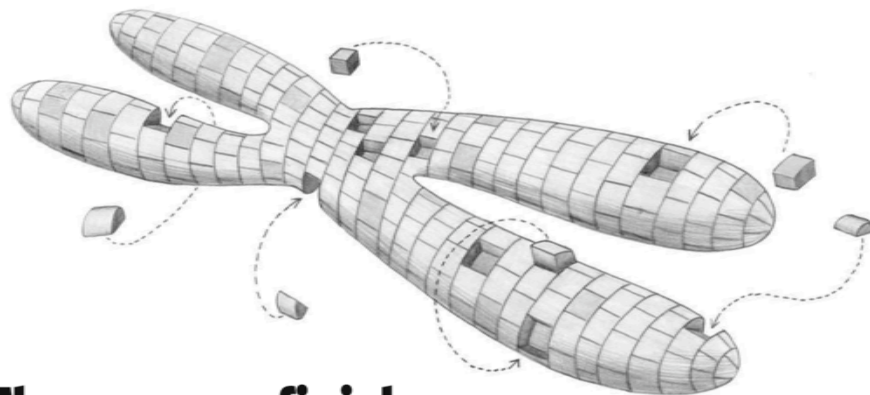
Finishing process

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 proceeds simultaneously at multiple scales, ranging from single nucleotides to the integrity of whole chromosomes. The fundamental challenge in genomic regions that are not well represented or readily resolved through random shotgun sequencing tend to be highly enriched in problematic sequences. Finishing such regions required the development of special approaches, which evolved substantially over time and varied among centres.

Essentially the finishing process involved two distinct components: (1) producing finished maps, consisting of contigs and accurate maps of overlapping long-insert clones spanning the euchromatic region of each chromosome arms and (2) producing finished clones, consisting of contigs and accurate maps spanning 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 components are described in Boxes 1 and 2. Further information about the finishing process and finishing standards can be found in the Supplementary Information (Note 1) and at <http://www.genome.gov/090922>.

In total, we generated a shotgun sequence from 56,204 large-insert clones (total length ~5.14 gigabases (Gb)) and finished the sequence from 45,742 of those clones (total length ~3.07 Gb). The clones contained primarily ~100-fold enriched, critical chromosomes 4

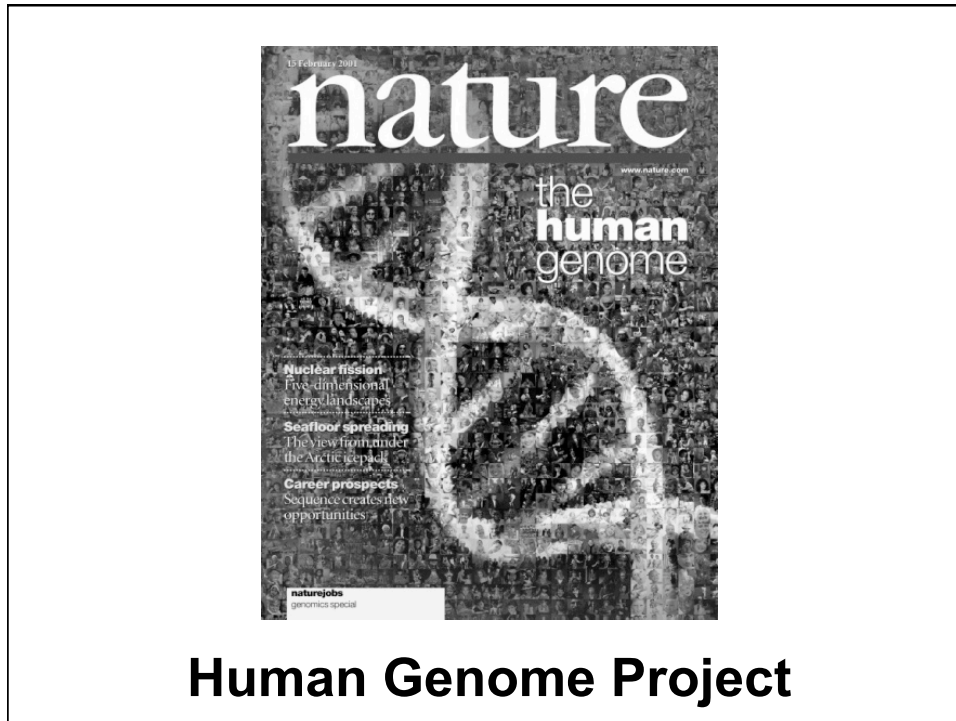
International Human Genome Sequencing Consortium (2004)



The genome finishers

Dedicated scientists are working hard to close the gaps, fix the errors and finally complete the human genome sequence. **Elie Dolgin** looks at how close they are.

Nature (2009)



The Path to Genomic Medicine



HGP



Realization of
Genomic Medicine

Genomic Medicine

***Healthcare tailored to the individual
based on genomic information***



The Path to Genomic Medicine



HGP



Realization of Genomic Medicine



“Fulfilling the Promise”



**Mapping
the Human Genome**

~1990 to ~2000

**Sequencing
the Human Genome**

~1998 to ~2003

**Interpreting
the Human Genome
Sequence**

~2003 to ???

The Human
Genome Project

Beyond
The Human
Genome Project

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

```
TGCCCGGAACTTTTCGGCTCTAAGGCTGATTTTGTATATACGAAAGGCACATTTTCCCTCCCTTTCAAATGCACCTTGCAACGTAAACAG
GAACCCGACTAGGATCATCGGAAAAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGG
CGCGAAGGAGGGTCTAGGAAGCTCTCCGGGAGCCGGTTCTCCCGCGGTGGCTTCTGTCTCCAGCGTTGCCAAGCTGAACTAAAGAGAGG
CCGCGACTGTCCGCCACCTGCGGGATGGCCTGGTCTGGCGGTAAGGACACGGACCTGGAAGGAGCGCGCGGAGGAGGGCTGGGAGTC
AGAATCGGAAAAGGAGGTGCGGGCGCGGAGGGAGCGAAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGG
GAAAGCCCTAGAGCAAAATTTGGGGCCGACAGGACAGCTCGGCTTTTAACTGGGAGTGAAGCGGGGAAAAGCAAAAGGAAAGGGGTGG
TGTGCGGAGTAGGGTGGTGGGGGAATTGGAAGCAAAATGACATCACAGCAGGTCAGAGAAAAGGGTTGAGCGGACAGCCACAGAGTAGTAG
GTCTTTGCAATAGGAGCTGAGCCAGACGGCCTAGCAGGACCCAGCGCCGAGAGCCATGCAGAGGTCCGCTCTGAAAAGGCCAGCGT
TGCTCCAAACTTTTTTCAGGTGAGAGGTGGCCACCAGGCTTCGGAAGACACGTCGCCACGAAAGAGAGGGCGTGTATGGGTGGGT
TGGGTAAAGGAAAGCAAGTTTTTAAAGATGCGCTATCATTCATTTTGGAAAGAAATGGGTATTTAGAAATAAAACAGAAAGCAATTA
AGAAGAGATGGAAGAATGAACTGAACTGATTAATAGAGAGCCATCTACTTGCACCTGAAAAGTTAGAACTCAAGACTCAAGTACGCTACT
ATGCACTGTTTTTATTTCATTTTTCTAAGAACTAAAATACTTGAATAAGTACCTAAGATATGGTTTATGGTTTTCCCCCTTCATGCCTGG
ACACTTGATTTCTTGGCACATACAGGTGCAATGCCTGCATATAGTAAGTGTCTAGAAAACATTTCTTGACTGAATCAGCCAAACAAAATTT
TTGGGTAGGTAGAAAATATATGCTTAAAGTATTATTGTTATGAGACTGGATATATCTAGTATTTGTCACAGGTAATGATTTCTCAAATTTG
AAAGCAAAATTTGTTGAAATATTTTAAAGAAAGTTACTTCACAAGCTATAAATTTTAAAGCCATAGGAATAGATACCGAAGTTATATCCAA
CTGACATTTAATAAATTTGATTCATAGCTAAATGTGATGAGCCACAGAAGCTTGCACAACTTAAATGAGATTTTTAAAATAGCATCTAAGTTCGG
AATCTTAGGCAAGTGTGTTAGATAGCACTCATATTTGAAGTGTCTTTGGATATGCATCTACTTTGTTCCCTGTTATTATACGTGTGTA
ATGAATGAATAGTACTCTCTCTTTGGACATTAATTGACACATAATACCCAATGAATAAGCATACTGAGGTATCAAAAAGTCAAAATATGT
TATAAATAGCTCATATATGTGTGAGGGGGAAGGAATTTAGCTTTCACATCTCTCTTATGTTAGTCTCTGCATGTGCGAGTTAATCTGGAAC
TCCGGTGTAAAGGAGAGACTTTGGCCCTGGAAGGAGACTCCCTCCCTGGATGAGAGAGAAGGACTTTACTTTTGGAAATATCTTTTTGTTG
TGATGTTATCCACCTTTGTTACTCCACTATAAATCGGCTTATCTATGATCTGTTTTCCTAGTCTTATAAAGTCAAAATGTTAATGGCAT
AAATATAGACTTTTTAGCAGAGAACTTTGAGAACTAAATGCCAACCAGTCTAAAATGCAGTTTTCAGAAAGTAAATATTTCACTGGATA
GTTCTAAATCTAATGAACTTAAAATAGCTTACTATGATCTGTCAAAGTGGTTTTTATAAATTTCTTTTACAAATCACCTGACACATTT
AATATAGGTTAAAATGCTATCAGCTGGTTGCAAGAAAATGATTACAAGGCTGCTAAGTGTGTTAAGAGCATACTCATTTCTGTTCTCC
AAAATATTCATAAGTGTCTTTAAGAATAGGATGTTTTTAAAGTAAAGTCTTACTATTTATAGGAAGTACCAATCACCTAAAATACCAATGA
TTCAAACTCCCTTGGCCCTTGGACTGCAATCTAAAAGTGTAAAACATATTTCTGCATTAAGTTAGGCAATGATGCTTAGTTTCAAA
GTGGTAGCTTTGGAGTCAGATTTTTGATTCAGATCTACATCTACTGTTAGTAGCTCTGTGCCCTGAGGACAGGTCCTTAAACATCTCTGTG
TGTACTTGCCTTTAAAATTTGGAGACTGTCATAGGGGTTAATCCCTTGAGAAAATGAATGTGAAAAGTTAGCCTAATGTTAATCTGCTATT
ATGGATACCATATTTTCACATTCACAGTACATGCACCTTGTTAATAATAGATGCTCAATTCATCTTTGAGTATAATTTTGTGACTCTCAAT
CTGGATATGCAATGAGTGGCCGTATGAGAAATTAATTTTAAAGAAATGTTGTTTTACATGGCTTACCAGATATACAGGAAACAGCTCACATG
TTTTATGATGTTTAAATGCTTAGAATTTAACTTTCTGAATAGGATCCCTTCAGTTTGGAGTCAAAAAGAGTAAAATTTATGTTAT
```

The Human Genome... by the Numbers

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

5% of 3B Bases = ~150M Bases

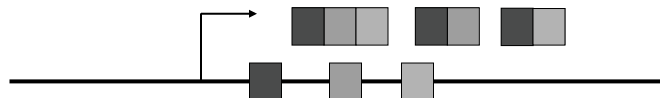
Do NOT Yet Know the Position of these ~150M Functional Bases

Lower Bound for the Amount that is Functional

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

```
TGCCCGGAACTTTTCGGCTCTAAAGGCTGATTTTTGATATACGAAAGGCACATTTCCCTTCCCTTTCAAATGCACCTTGCAAACGTAAACAG
GAACCCGACTAGGATCATCGGAAAAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGG
CGCGAAGGAGGGTTAGGAAGCTCTCCGGGAGCCGGTTCTCCCGCGGTGGCTTCTCTGCTCCAGCGTGGCCAACTGAACCTAAAGAGAGG
CCCGCACTGTCCGCCACTGCGGGATGGCCCTGGTCTGGGCGGTAAGGACACGGACCTGGAAGGAGCGCCGCCAGGAGGAGGCTGGGAGTC
AGAATCGGAAAAGGAGGTGCGGGGCGGCAGGGAGCGAAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGG
GAAAGCCCTAGAGCAAATTTGGGCGCGACAGGCAGCACTCGCTTTTAACTGGGCACTGAAGGCGGGGAAAAGAGCAAAAAGGAAAGGGGTTGG
TGTGCGGAGTAGGGTGGTGGGGGAATTGGAAGCAAATGACATCACAGCAGGTCAGAGAAAAAGGTTGAGCGGCAGGCACCCAGAGTAGTAG
GTCTTTGGCATTAGGAGCTGAGCCGAGCGGCCCTAGCAGGACCCAGCGCCGAGAGACCATGCAGAGGTCGCTCTGAAAAGGCCAGCGT
TGCTCCAAACTTTTTTCAGGTGAGAGGTGGCAACCGAGCTTCCGAAAGACACGTGCCACGAAAAGAGGAGGCGCTGTGATGGGTTGGGTT
TGGGTAAAGGAATAAGCAGTTTTTAAAGATGCGCTATCATTCATTGTTTTGAAAAGAAAATGGGGTATTGTAGAATAAAAACAGAAAGCATTAA
AGAAGAGATGGAAGAATGAAGCTGAAGTGAATGAAATAGAGAGCCATCTACTTGCACCTGAAAAGTTAGAACTCAAGACTCAAGTACGCTACT
ATGCACTGTTTTATTTCATTTTTCTAAGAACTAAAAATCTTGAATAAGTACCCTAAGATGTTTTATGGTTTTCCCCCTTCATGCTTGTG
ACACTTGATTTCTTTCGCGACATACAGGTGCCATGCTGCATATAGTAAAGTGTCTAGAAAACATTTCTTGACTGAATCAGCCAAACAAAATTT
TTGGGTAGGTAGAAAATATGCTTAAAGTATTATTGTTATGAGACTGGATATATCTAGTATTTGCACAGTAAATGATCTCTCAAAAATTTG
AAAGCAAATTTGTTGAAATATTTTAAAGAAAGTTACTTCACAGCTATAAAATTTAAAGCCATAGGAATAGATACCGAAGTTATATCCAA
CTGACATTTAATAAAATGATTCATAGCTAAATGATGAGCCACAGAGCTTGCACAACTTAAATGAGATTTTTAAAATAGATCAATAGGTTGCG
AATCTTAGGCAAGTGTGTTAGATGAGCACTTCATATTTGAAAGTGTCTTTGATATGCACTACTTTGTTCCCTGTTATTATACTGGTGTGA
ATGAATGAATAGTACTCTCTCTCTGGACATTACTTACACATAATACCCTAATGAATAGCATACTGAGGATCAAAAAGTCAAAATATGTT
TATAAATAGCTCATATATGTGTAGGGGGGAAGGAATTTAGCTTTCACATCTCTCTTATGTTTGTCTCTGCATGTCAGTTAATCCTGGAAC
TCCGGTGCTAAGGAGAGACTGTGGCCCTTGAAAGAGAGCTCCTCCCTGGGATGAGAGAGAAGGACTTACTCTTGGAAATTCATTTTTGTTGTG
TGATGTTATCCACCTTTGTTACTCCACTATAAAATCGGCTTATCTATGATCTGTTTTCTTAGTCTTATAAAGTCAAAAATGTTAATGGCATT
AAATATAGACTTTTTAGCAGAGAACTTTGAGGAACCTAAATGCCAACCAGTCTAAAATGCAGTTTTGAGAGAAATGAATATTTGATGGATA
GTTCTAATACTAATGAACTTAAAATAGCTTACTATGATCTGTCAAAGTGGGTTTTATATAATTTCTTTTTACAAATCACCTGCACACATTT
AAATAGGTTAAAAGTCTATCAGCTGGTTTTGCAAGAAAATGATATCAAAAGGCTGCTAAGTGTGTTAAGACATACTCATTCTGTTCTCC
AAAATATTTCTAAGGTGTTTTAAGAATAGGTATGTTTTAAAGTTAAGTTCCTACTATTATAGGAAGTACCAATCACCTAAAATACCAATGA
TTCAAACTTCCCTTGGCCCTTGGACTGCAATTTCAAAAGTGTAAAACATATAATTTCTGCATTAAAGTAGGCACTATGCTTAGTTTTCAA
GTGGTAGGCTTGGAGTCAGATTTTTGATTCAGATCTACATCTACTGTTAGTAGCTCTGTTGCTGAGGAGGTCCCTTAACATCTCTGTG
TGTGACTTGACCTTAAAATTTGAGAGCTGTCATAGGGTAAATCCCTTGAGAAAATGAATGTGAAAAGTTAGCCTAATGTTAATCTGCTATTAT
ATGATTACCATATTTTACATTCACAGTACATGCACTTGGTAAATAAGATGCTCAATTCATCTTTGAGTATAATTTGCTGACTCTCAAT
CTGGATATGCAATGAGTGGCCTGTATGAGAAATTAATTTATGAAAATTTGTTTTCACATGGCCCTTACCAGATATACAGGAAACACGTCACATG
TTCTTATGTTATGTTAAATGCCTTAGAAATTAACTTTCTGAAATAGGATCCTTCACTTGTAGAGTCATAAAGAGTAAAATTAATGTTAT
```

Coding Sequences (i.e., Genes)



		Second Letter				
		T	C	A	G	
First Letter	T	TTT Phe TTC Leu TTA Leu TTG Leu	TCT Ser TCC Ser TCA Ser TCG Ser	TAT Tyr TAC Stop TAA Stop TAG Stop	TGT Cys TGC Stop TGA Stop TGG Trp	T C A G
	C	CTT Leu CTC Leu CTA Leu CTG Leu	CCT Pro CCC Pro CCA Pro CCG Pro	CAT His CAC Gln CAA Gln CAG Gln	CGT Arg CGC Arg CGA Arg CGG Arg	T C A G
	A	ATT Ile ATC Met ATA Met ATG Met	ACT Thr ACC Thr ACA Thr ACG Thr	AAT Asn AAC Lys AAA Lys AAG Lys	AGT Ser AGC Arg AGA Arg AGG Arg	T C A G
	G	GTT Val GTC Val GTA Val GTG Val	GCT Ala GCC Ala GCA Ala GCG Ala	GAT Asp GAC Glu GAA Glu GAG Glu	GGT Gly GGC Gly GGA Gly GGG Gly	T C A G

The Genetic Code

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

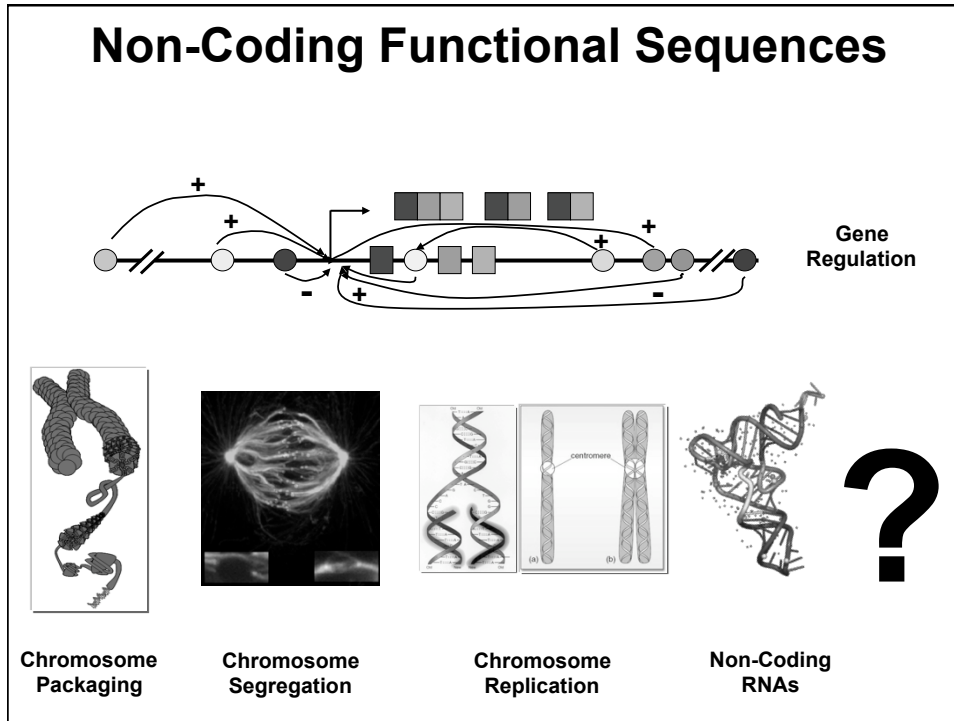
Good Inventory at Present

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

```

TGCCGCGGAACCTTTTCGGCTCTAAGGCTGATTTTGTATATACGAAAGGCACATTTTCCTTCCCTTTTCAAATGCACCTTGCAACGTAACAG
GAACCCGACTAGGATCATCGGAAAAGGAGGAGGAGGAAGGCAGGCTCCGGGAAGCTGGTGGCAGCGGGTCTGGGTCTGGCGGACCCCTGA
CGCGAAGGAGGGTCTAGGAAGCTCTCCGGGAGCCGGTTCFCCCGCGGTGGCTTCTGTCTCCAGCGTTGCCAAGTAAAGAGAGG
CCGCGACTGTCGCCACCTGCGGGATGGCCCTGGTGTGGCGGTAAGGACACGGACCTGGAAGGAGCGCGCGGAGGGAGGGAGGCTGGGAGTC
AGAATCGGAAAGGAGGTCGCGGGCGCGGAGGAGCGAAGGAGGAGGAGGAGGAGCGGGAGGGCTGCTGGCGGGGTGCGTAGTGGGTGGA
GAAAGCCGCTAGAGCAAAATTTGGGGCCGACCGAGCAGCAGCTCGCTTTTAACTGGGCACTGAAGCGGGGAAAGAGCAAAAGGAAGGGGTGG
TGTGCGGAGTAGGGTGGTGGGGGAATTTGGAAGCAATGACATCACAGCAGGTCAGAGAAAAGGGTTGAGCGGCAGGCACCCAGAGTAGTAG
GTCTTTGGCATTAGGAGCTTAGGCCAGACGGCCCTAGCAGGGACCCAGCGCCGAGAGACCATGCAGAGGTCGCTCTGAAAAGGCCAGCGT
TGTCTCAAACCTTTTTTCAGGTGAGAAGGTGGCCAACCGAGCTTCGGAAGACACGTGCCACGAAAGAGGAGGGCGTGTGTATGGTTGGGT
TGGGTAAAGGAATAAGCAGTTTTTAAAAGATGCGCTATCATTTGTTTTGAAAGAAAATGTGGTATTTAGTAATAAACAGAAAGCATT
AGAAGAGATGGAAGAATAACTGAAGCTGATTTGAATAGAGAGCCACTACTTGAACCTGAAAAGTTAGAATCTCAAGACTCAAGTACGCTACT
ATGCACTTGTTTTTATTTCATTTTCTAAGAAAATAAAAATACTTGTTAATAAGTACCTAAGTATGGTTATTTGGTTTTCCCTTCATGCCTTG
ACACTTGAATGTCTTCTGGCACATACAGGTGCCATGCCATAGTAAAGTGCFCAGAAAACATTTCTTGACTGAATTCAGCCAAACAAAAT
TTGGGTAGGTAGAAAATATATGCTTAAAGTATTTATTTGTTATGAGACTGGATATATCTAGTATTTGTCACAGGTAATGATTTCTCAAATAATG
AAAGCAAAATTTGTTAAATATTTATTTGAAAAGATTTACTTCAACAGCTATAAATTTTAAAAGCCATAGGAATAGATACCGAAGTTATATCCAA
CTGACATTTAATAAATTTGATTCATAGCCTAATGTGATGAGCCACAGAAGCTTGCAAACTTTAATGAGATTTTTTAAAATAGCATCTAAGTTCGG
AATCTTAGGCAAGTGTGTTAGATGTAGCACTTCATATTTGAAGTGTCTTTGGATATGTCATCTACTTTGTTCCCTGTATATATCTGGTGTGA
ATGAATGAATAGTACTGCTCTCTTGGGACATTTACTTGACACATAAATACCCAATGAATAAGCATACTGAGGTATCAAAAAGTCAAAATATGT
TATAATAGCTCATATATGTGTGATGGGGGAAGGAATTTAGCTTTCACATCTCTTATGTTAGTTCTCTGATGTGCAAGTAACTCCTGGAAC
TCCGCTGCTAAGGAGAGACTGTGGCCCTTGAAGGAGAGCTCCCTCCCTGTGGATGAGAGAGAGGACTTTACTCTTTGGAAATATCTTTTGTGT
TGATGTTATCCACCTTTTGTACTCAACTATAAAATCGGCTTATCTATTGATCTGTTTTCCCTAGTCTTATAAAAGTCAAAATGTTAATGGCAT
AAATATAGACTTTTTTAGCAGAGAATTTGAGAACTAATGCCAACGAGCTAAAAATGCAGTTTTTCAGAAAGATGAATATTTCTATGGATA
GTCTAAATACTAATGAACTTAAAATAGCTTACTATTGATCTGTCAAAGTGGTTTTTATATAAATTTCTTTTACAAATCACCTGACACATTT
AATATAGTTAAAAATGCTATCAGGCTGGTTTGAAGAAAATGTATTACAAAGGCTGCTAAGTGTGTAAAGAGCATACTCATTTCTGCTCC
AAAAATTTTCATAGGCTTTAAGAAATAGGATGTTTTTAAAAGTAAAGTCCCTACTATTTATAGGAACGACAACTCACCTAAAATACCAATGA
TTACAACCTCTCTGCGCTTCTGACTGCAATTTCAAAGTGTAAAAAACATATTTCTGCATTAAAGTAGGCAGTATTGCTTAGTTTTCAA
GTGGTAGGCTTTGGAGTCAGATTTTGTATTGATCTCAGATCTACTGTTTAGTAGCTCTGTGCGCTGAGGACGGTCCCTTAACATCTCTGTG
TGTGACTGACCTTTAAAATTTGGAGACTGTACATAGGGGTTAATCCCTTGAGAAAATGAATGTGAAAAGTTAGCCCTAATGTTAATCTGATATAT
AATAGTATACCAATTTTCAATTCACAGTACATGACCTTGTAAATATAAGATGCTCAATTCATCTTTAGATATAATTTTGTGACTCTCAAT
CTGGATATGCAATGAGTGGCCGTGATGAGAATTTAATTTATGAAAATTTGCTTTTACATGCGCTTACCAGATATACAGGAAACACCTCACATG
TTTTATGATGTTGTTAAATGCTTAGAATTTAATTTCTGAATAGGATCCCTCAGTTTGGAGTCAATAAAGAGTAAAATTTATGTTAT

```



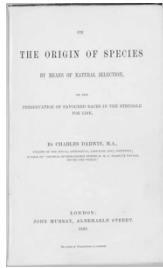
The Human Genome... by the Numbers

~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



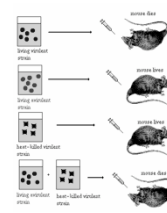
Mendel

1865



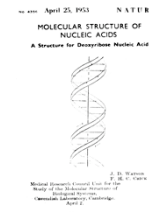
Miescher

1871



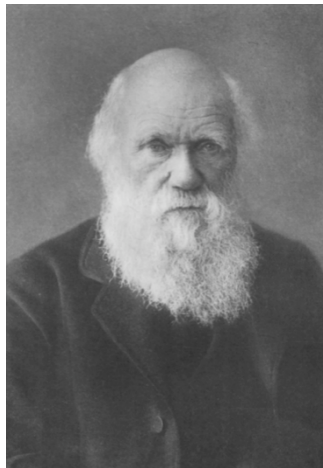
Avery

1944



Watson & Crick

1953

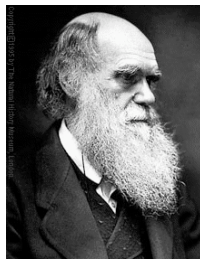


Charles Darwin
Born February 12, 1809



"It is not the strongest of the species that survives, nor the most intelligent that survives. It is the one that is the most adaptable to change."

(Attributed to Darwin)

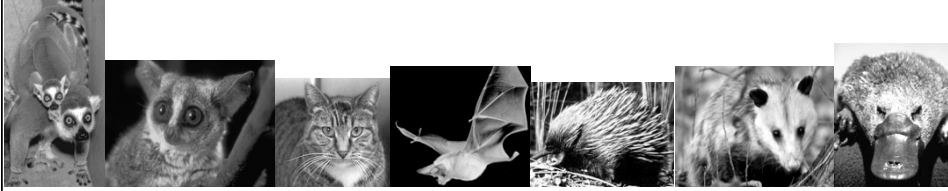


Charles Darwin (1809-1882)

"For the last three and a half billion years, evolution has been taking notes."

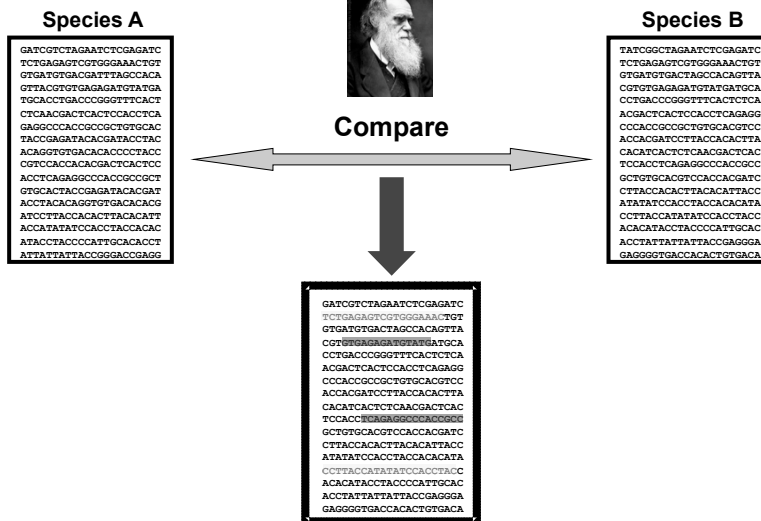
—Eric Lander

Inter-Species Sequence Comparisons



Comparative Sequence Analysis

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

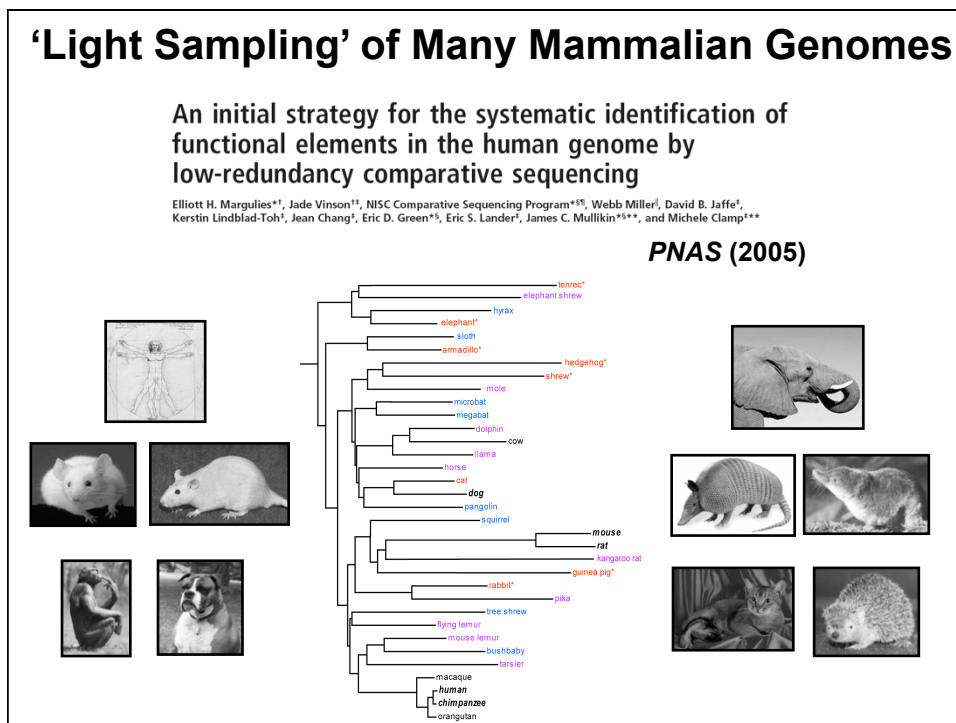
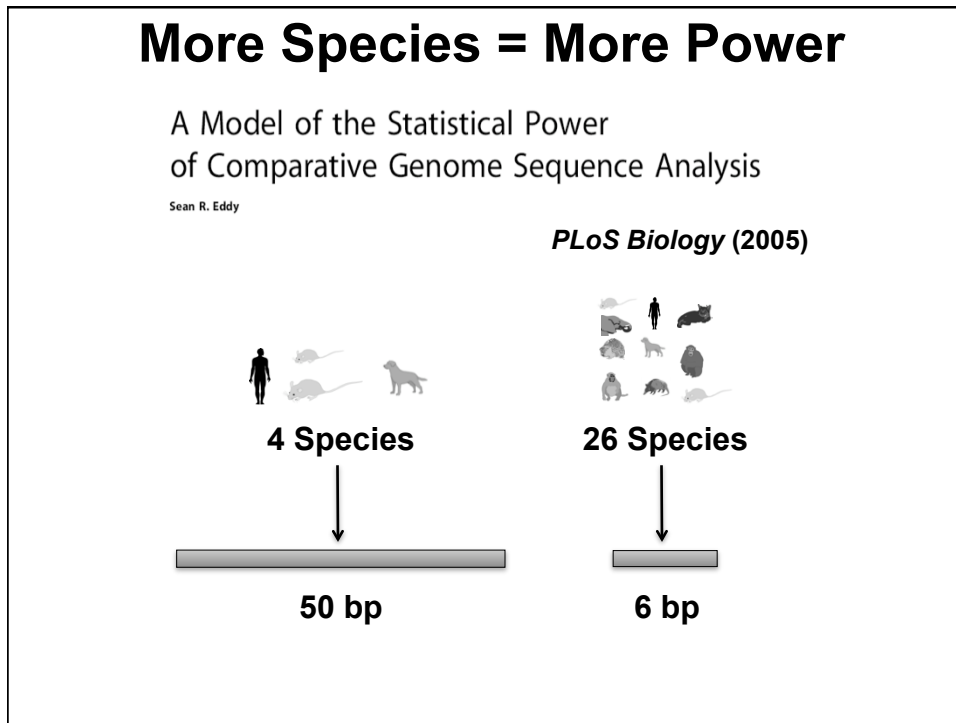


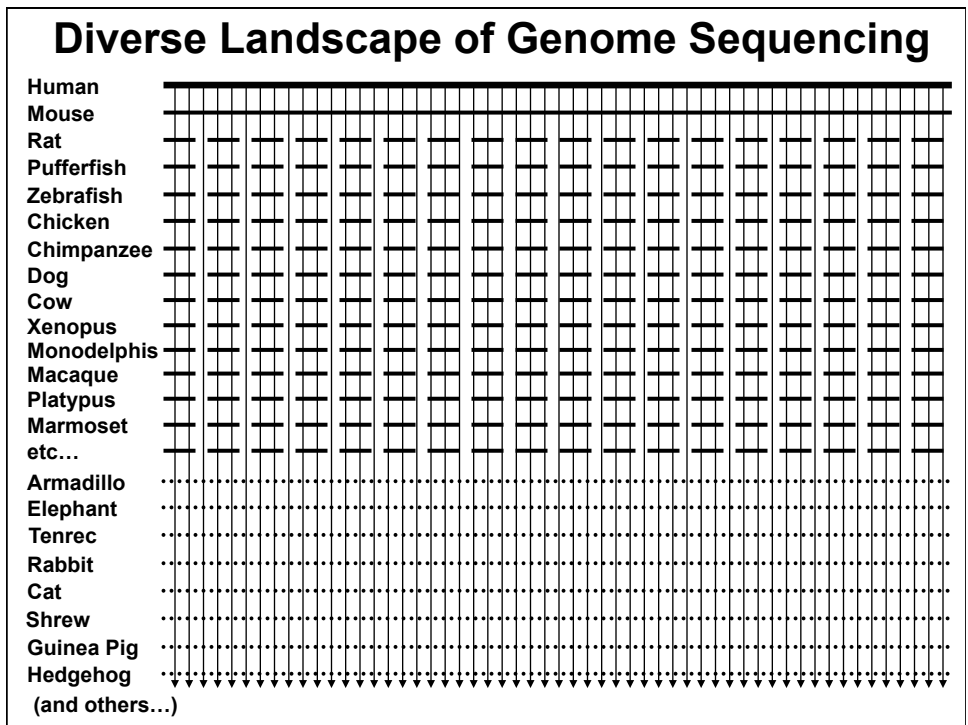
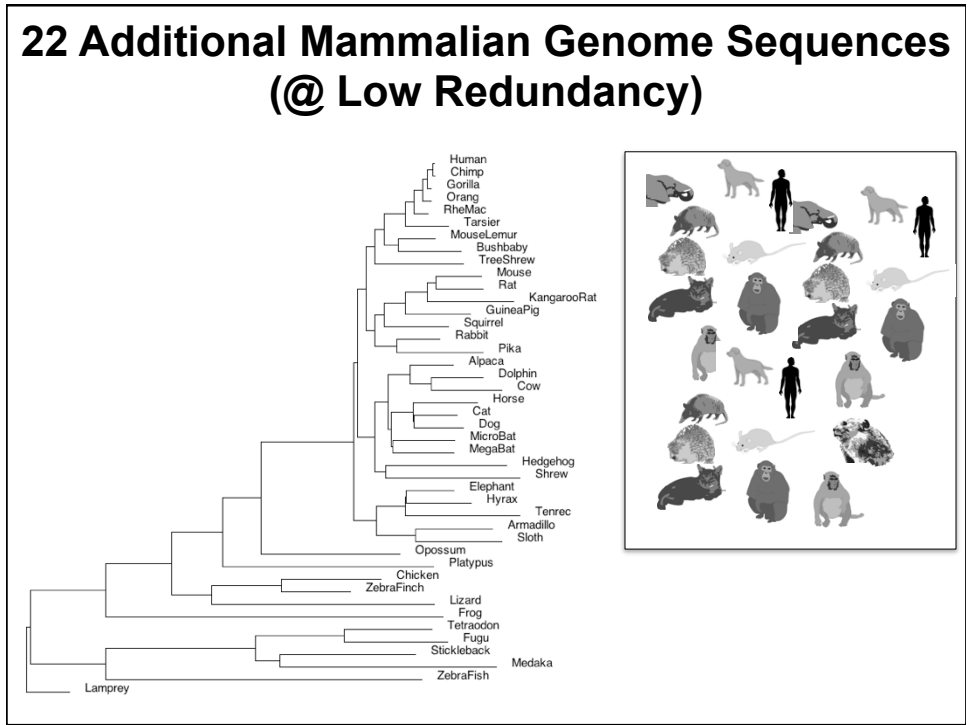
Vertebrate Genome Sequences



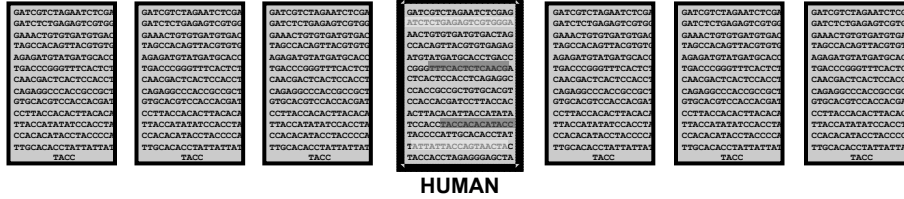
Diverse Landscape of Genome Sequencing

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

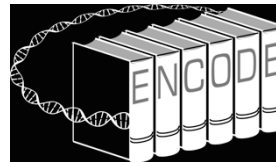




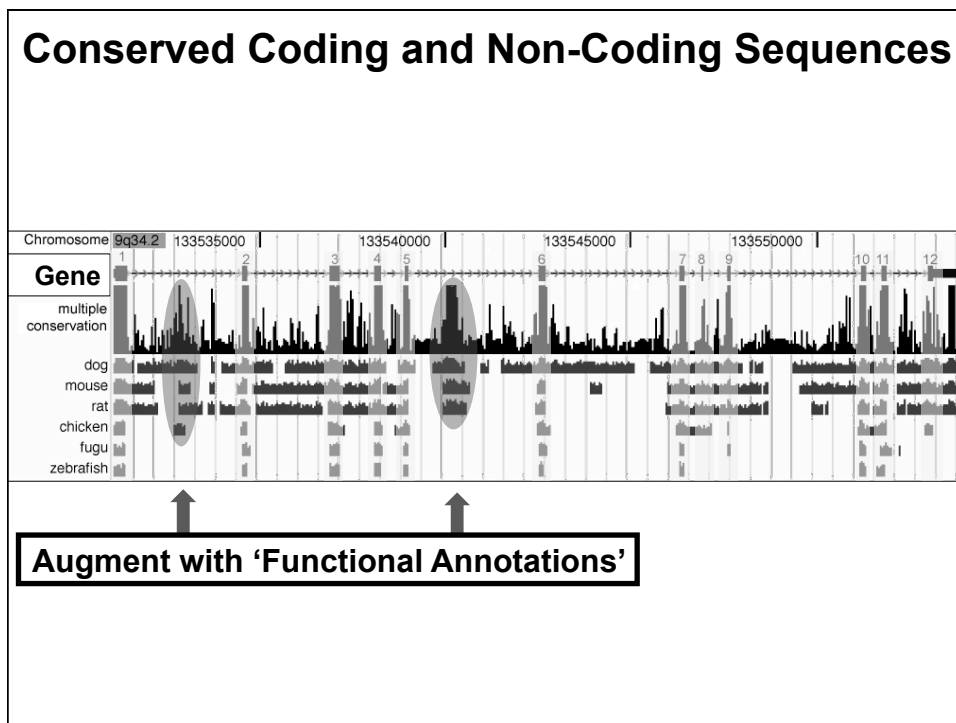
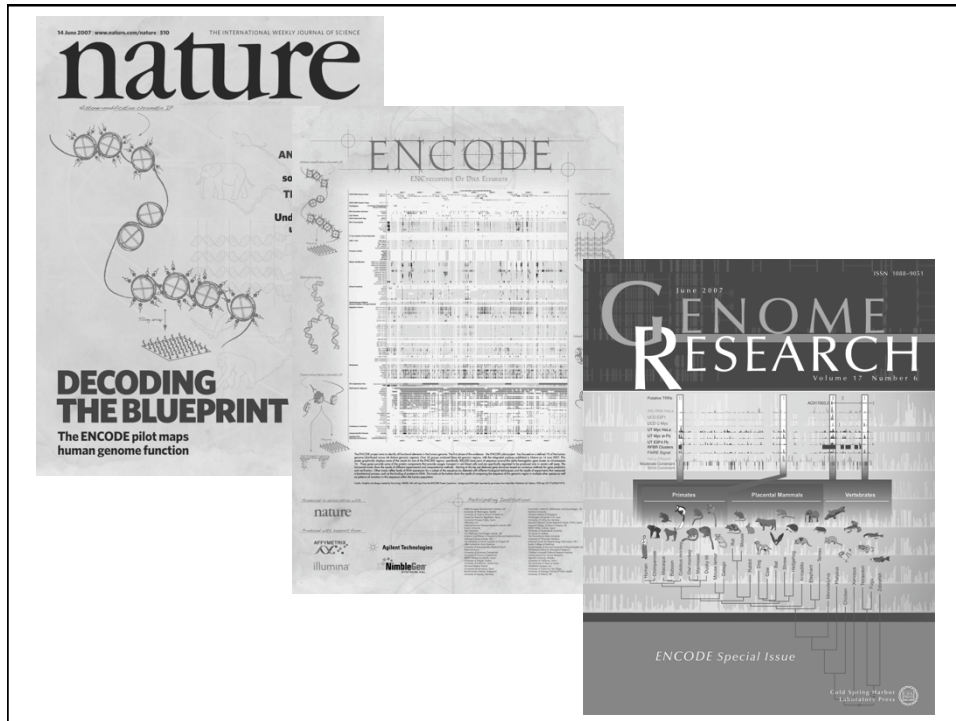
Multi-Species Sequence Comparisons



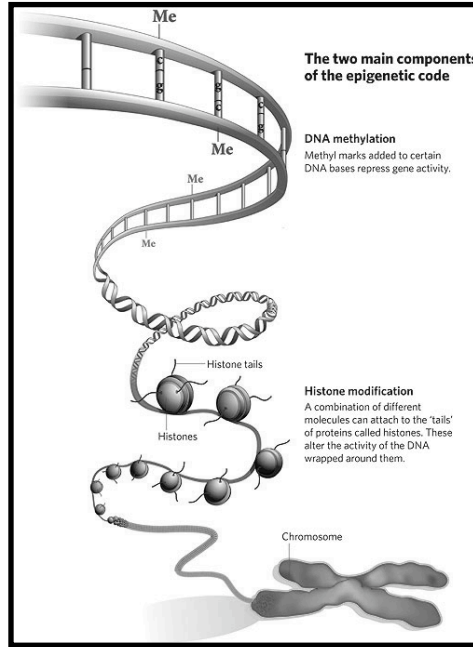
ENCODE Project



- 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



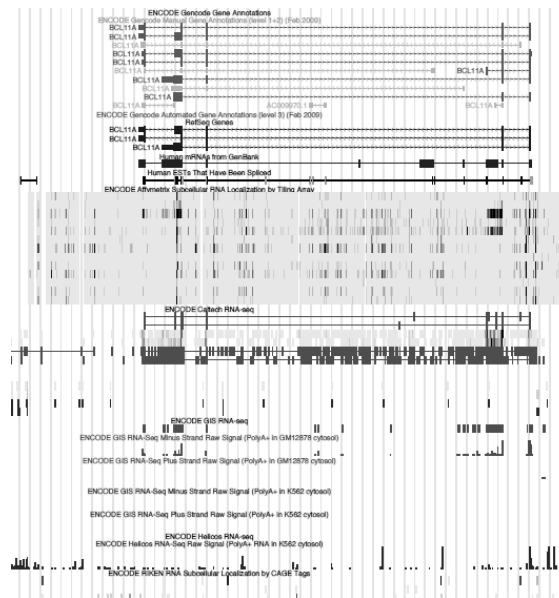
The Epigenetic Landscape

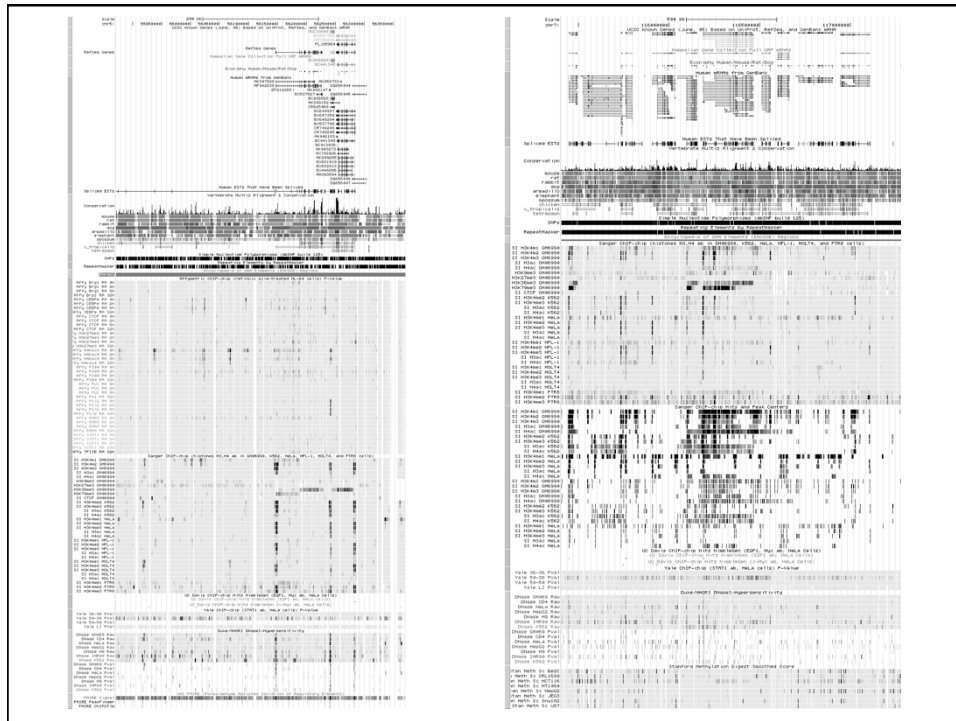


ENCODE: Lots of Data and Data Types

Generated by:

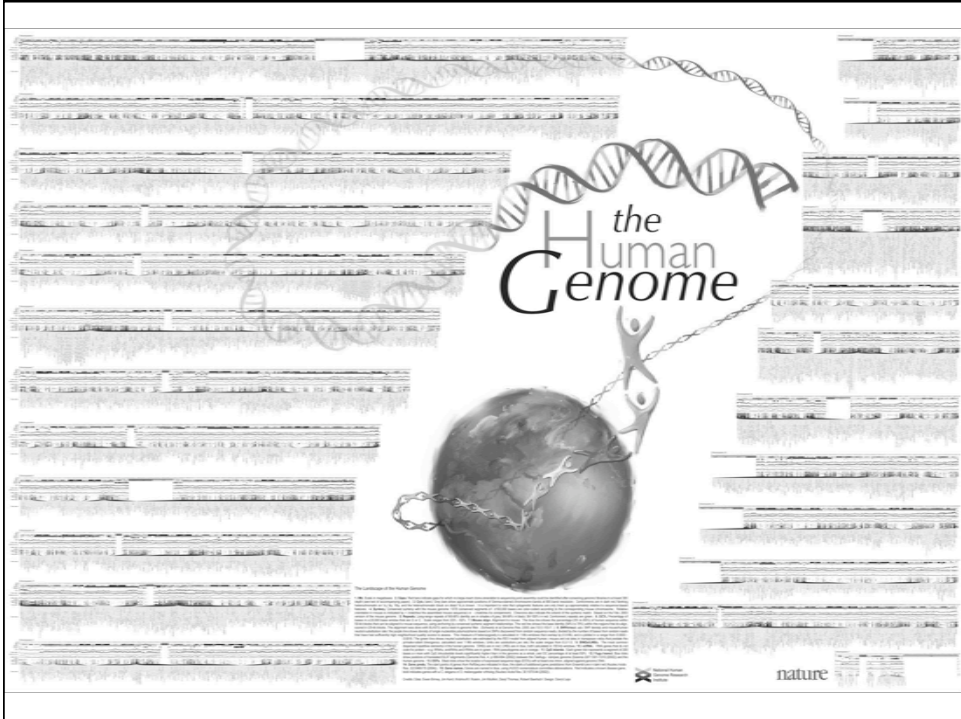
- RNA-Seq
- RNA-array
- TF ChIP-Seq
- Histone modif ChIP-Seq
- DNaseHS-Seq
- FAIRE-Seq
- Methyl-Seq
- Methyl27-bisulfite
- etc.
- etc.
- etc.



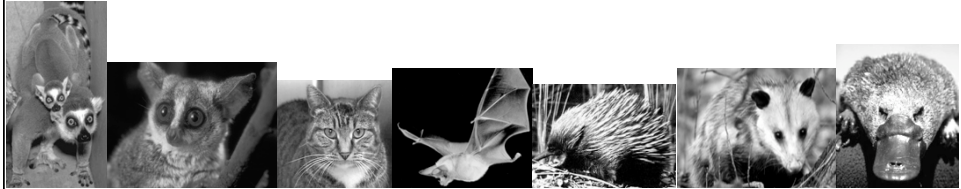


Expanding ENCODE Portfolio





Inter-Species Sequence Comparisons



The Genomics of Human Evolution



(David Haussler, Stephen O'Brien, & Oliver Ryder)

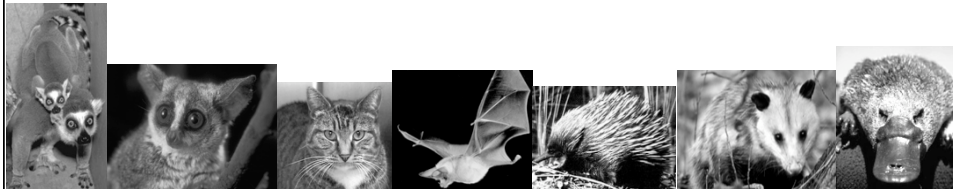


Genome 10K: A Proposal to Obtain Whole-Genome Sequence for 10 000 Vertebrate Species

GENOME 10K COMMUNITY OF SCIENTISTS*

J. Heredity (2009)

Inter-Species Sequence Comparisons



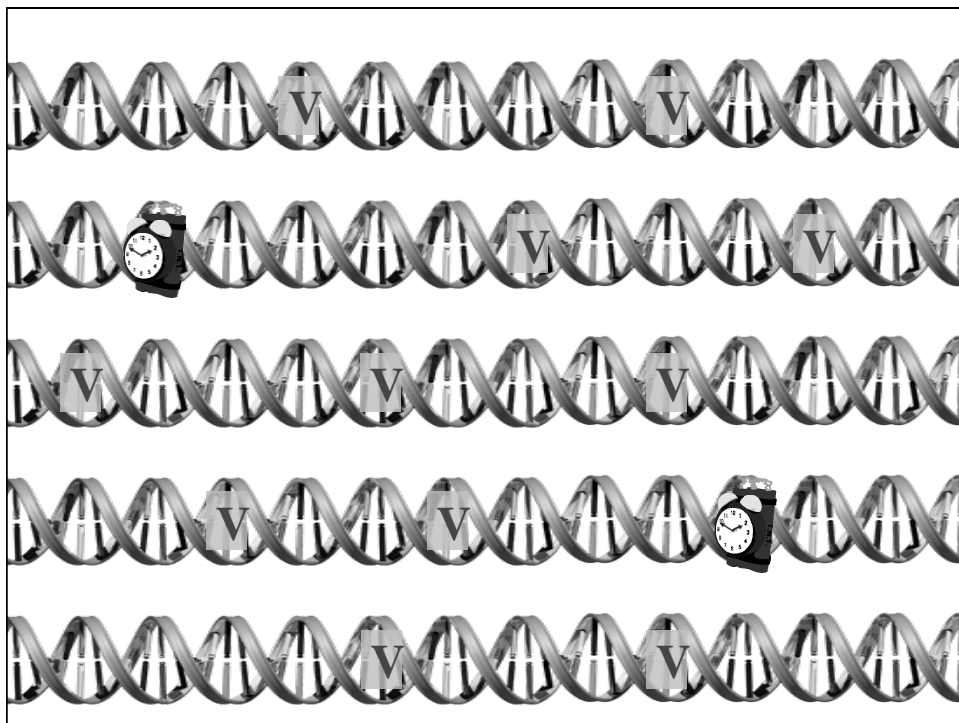
Intra-Species Sequence Comparisons



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



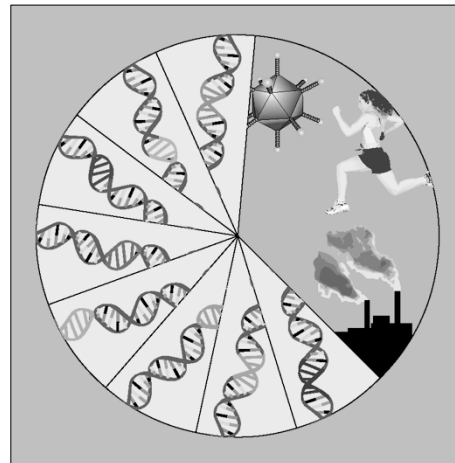
all of us carry a significant number of 'glitches' in our genomes.



Genomic Architecture of Genetic Diseases

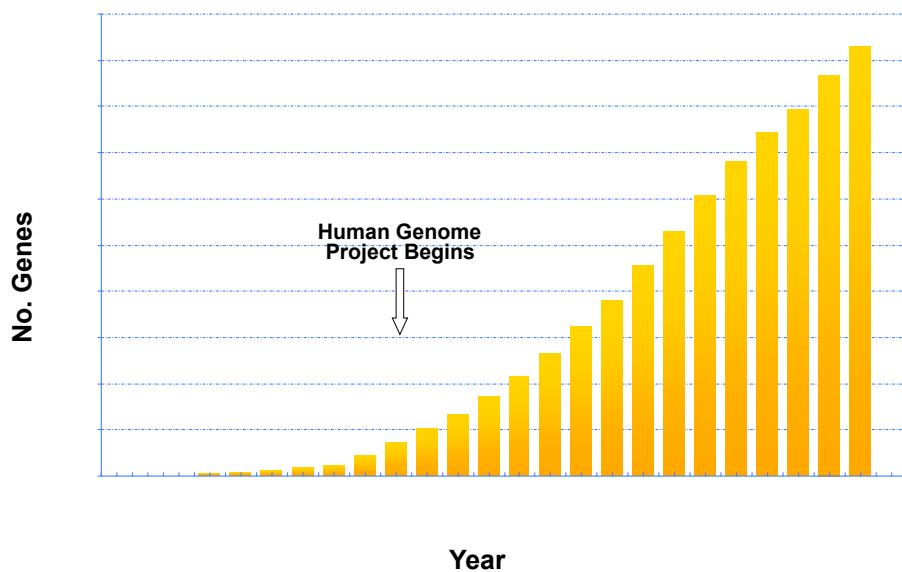


Rare, Simple, Monogenic, Mendelian...



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

Human Disease Genes Identified: 1981-2005



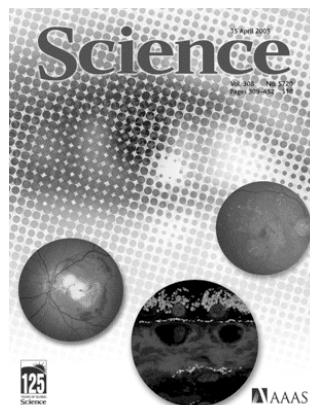
Source: Online Mendelian Inheritance in Man

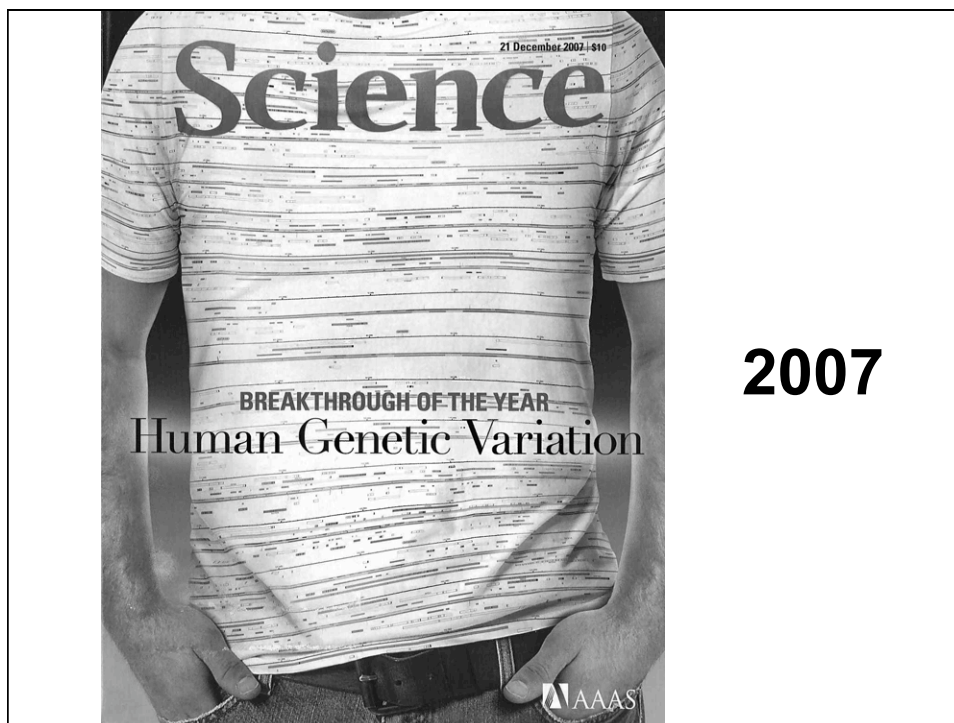
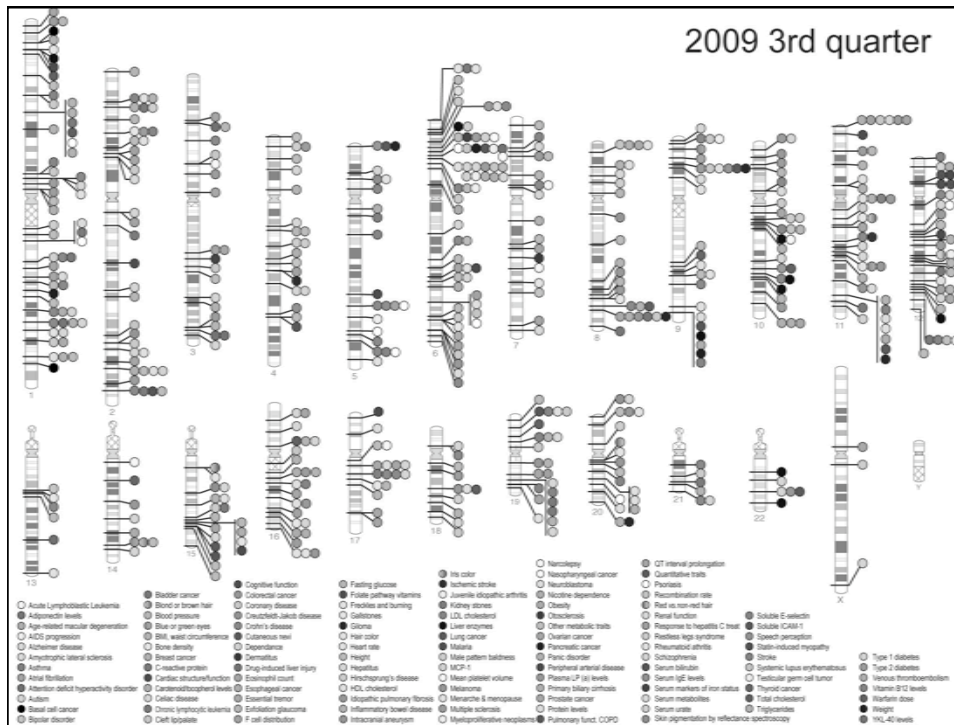


The First HapMap Success Story: Age-Related Macular Degeneration

Complement Factor H Polymorphism in Age-Related Macular Degeneration

Robert J. Klein,¹ Caroline Zeiss,^{2*} Emily Y. Chew,^{3*} Jen-Yue Tsai,^{4*} Richard S. Sackler,¹ Chad Haynes,¹ Alice K. Henning,⁵ John Paul SanGiovanni,³ Shrikant M. Mane,⁶ Susan T. Mayne,⁷ Michael B. Bracken,⁷ Frederick L. Ferris,³ Jurg Ott,¹ Colin Barnstable,² Josephine Hoh^{7†}



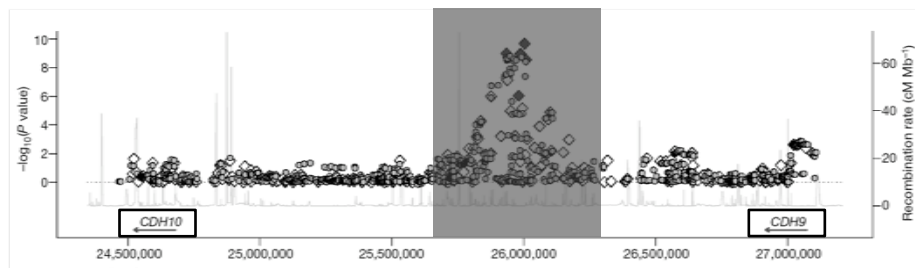


Genetic Association within Intergenic Region

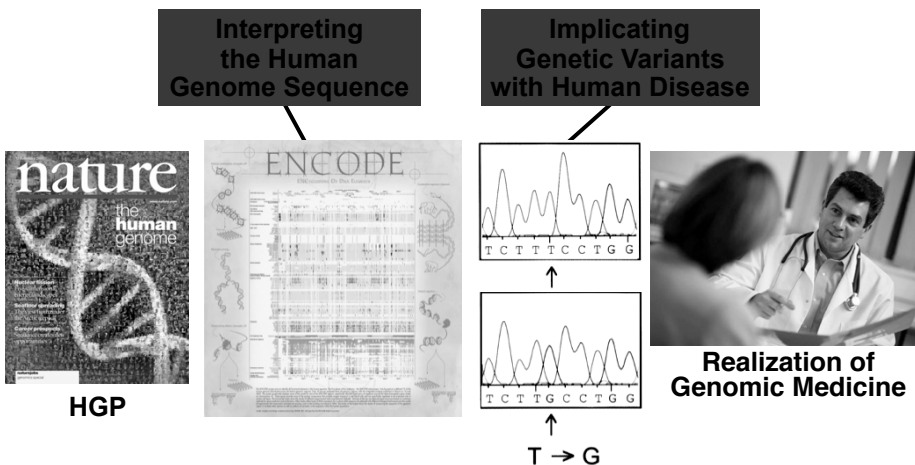
Common genetic variants on 5p14.1 associate with autism spectrum disorders

Kai Wang^{1*}, Haitao Zhang^{1*}, Deqiong Ma^{2*}, Maja Bucan³, Joseph T. Glessner¹, Brett S. Abrahams⁴, Daria Salyakina², Marcin Imielinski¹, Jonathan P. Bradfield¹, Patrick M. A. Sleiman¹, Cecilia E. Kim¹, Cuiping Hou¹, Edward Frackelton¹, Rosetta Chiavacci¹, Nagahide Takahashi⁵, Takeshi Sakurai⁵, Eric Rappaport⁶, Clara M. Lajonchere⁷, Jeffrey Munson⁸, Annette Estes⁹, Olena Korvatska⁸, Joseph Piven⁴, Lisa I. Sonnenblick⁴, Ana I. Alvarez Retuerto⁴, Edward I. Herman⁴, Hongmei Dong⁴, Ted Hutman⁴, Marian Sigman⁴, Sally Ozonoff¹⁰, Ami Klin¹¹, Thomas Owley¹², John A. Sweeney¹², Camille W. Brune¹², Rita M. Cantor¹³, Raphael Bernier⁸, John R. Gilbert², Michael L. Cuccaro², William M. McMahon¹⁴, Judith Miller¹⁴, Matthew W. State¹¹, Thomas H. Wassink¹⁵, Hilary Coon¹⁴, Susan E. Levy⁶, Robert T. Schultz⁶, John I. Nurnberger Jr.¹⁶, Jonathan L. Haines¹⁷, James S. Sutcliffe¹⁸, Edwin H. Cook¹², Nancy J. Minshew¹⁹, Joseph D. Buxbaum^{5,20}, Geraldine Dawson², Struan F. A. Grant^{1,6}, Daniel H. Geschwind⁴, Margaret A. Pericak-Vance², Gerard D. Schellenberg²¹ & Hakon Hakonarson^{1,6}

Nature (2009)




The Pathway to Genomic Medicine

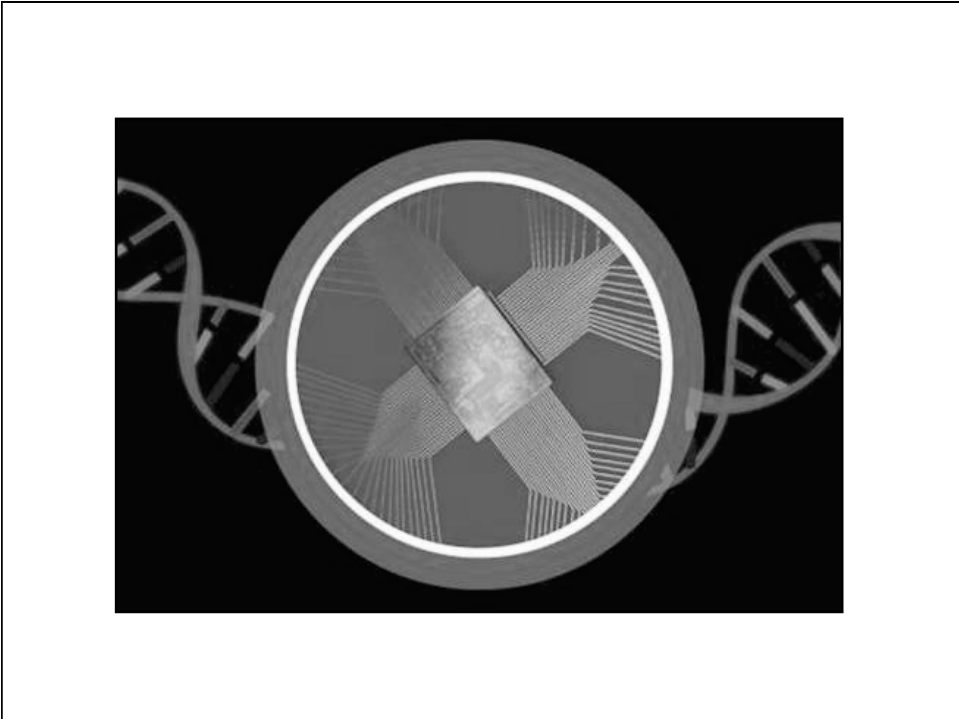
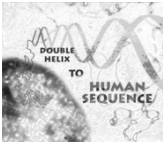


Human Genome Sequence

>\$1,000,000,000



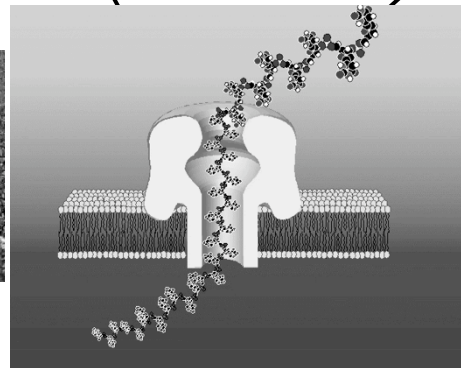
~\$1,000



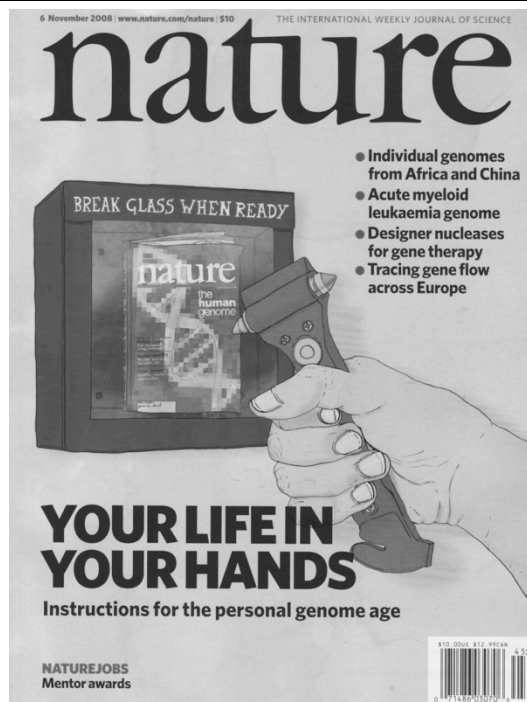
The Path to Genomic Medicine

Interpreting
the Human
Genome Sequence

Implicating
Genetic Variants
with Human Disease



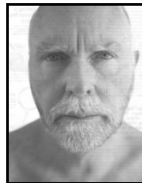
Realization of
Genomic Medicine



Individual Genome Sequences

The Diploid Genome Sequence of an Individual Human

Samuel Levy¹, Granger Sutton¹, Pauline C. Ng¹, Lars Feuk², Aaron L. Halpern³, Brian P. Walenz⁴, Nelson Axelrod¹, Jaal Husni⁵, Ewan F. Kirkness⁶, Genady Denisov⁷, Yuan Lin¹, Jeffrey R. MacDonald⁸, Andy Ming Chin Pang⁹, Mary Shago¹⁰, Timothy B. Stockwell¹¹, Alexia Tsamirani¹², Vinset Bafna¹³, Vikas Bansal¹⁴, Saul A. Kravitz¹⁵, Dana A. Busam¹⁶, Karen Y. Beeson¹⁷, Tina C. McIntosh¹⁸, Karin A. Remington¹⁹, Josef P. Abul²⁰, John Gill²¹, Jon Borman²², Yu-Hui Rogers²³, Marvin E. Frazier²⁴, Stephen W. Scherer²⁵, Robert L. Strausberg²⁶, J. Craig Venter²⁷



PLoS Biol (2007)

Accurate whole human genome sequencing using reversible terminator chemistry

A list of authors and their affiliations appears at the end of the paper

Nature (2008)

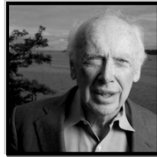
The diploid genome sequence of an Asian individual

Jun Wang^{1,2,3,4*}, Wei Wang^{1,2*}, Ruiqiang Li^{1,3,4*}, Yingrui Li^{1,3,4*}, Geng Tian^{1,2}, Laurie Goodman¹, Wei Fan¹, Junqing Zhang¹, Jun Li¹, Jianbin Zhang¹, Yiran Guo^{1,2}, Binxiao Feng¹, Heng Li^{1,2}, Yao Lu¹, Xiaodong Fang¹, Huiqing Liang¹, Zhenglin Du¹, Dong Li¹, Yang Zhao^{1,2}, Yujie Hu^{1,2}, Zhenchen Yang¹, Hancheng Zheng¹, Ines Hellmann¹, Michael Inouye¹, John Pool¹, Xin Yi^{1,2}, Jing Zhao¹, Jinjie Duan¹, Yan Zhou¹, Junjie Qin^{1,2}, Lijia Ma^{1,2}, Guoqing Li¹, Zhenhao Yang¹, Guojie Zhang^{1,2}, Bin Yang¹, Chang Yu¹, Fang Liang^{1,2}, Wenjie Li¹, Shaochuan Li¹, Dawei Li¹, Peisiang Ni¹, Jian Ruan¹, Qibin Li¹, Hongmei Zhu¹, Dongyuan Liu¹, Zhike Lu¹, Ning Li¹, Guangshu Gao^{1,2}, Jianguo Zhang¹, Jia Ye¹, Lin Fang¹, Qin Hao^{1,2}, Quan Chen^{1,2}, Yu Liang^{1,2}, Yeyang Su^{1,2}, A. san^{1,2}, Guo Ping^{1,2}, Shuang Yang¹, Fang Chen^{1,2}, Li Li¹, Ke Zhou¹, Hongkun Zheng¹, Yuanyuan Ren¹, Ling Yang¹, Yang Gao^{1,2}, Guohua Yang^{1,2}, Zhao Li¹, Xiaodi Feng¹, Karsten Kristiansen¹, Gane Ka-Shu Wong^{1,2}, Rasmus Nielsen¹, Richard Durbin¹, Lars Bolund^{1,11}, Xuqing Zhang^{1,6}, Songgang Li^{1,12}, Huaming Yang^{1,13} & Jian Wang^{1,14}

Nature (2008)

The complete genome of an individual by massively parallel DNA sequencing

David A. Wheeler^{1*}, Mithreyan Srinivasan^{2*}, Michael Egholm^{3*}, Yufeng Shen^{1*}, Lei Chen¹, Amy McGuire¹, Wen He¹, Yi-Ju Chen¹, Vinod Makhiani¹, G. Thomas Roth¹, Xavier Gomes¹, Karrie Tartaro¹, Fahem Nazif¹, Cynthia L. Turcotte¹, Gerard P. Irzyk¹, James R. Lupski^{1,4}, Craig Chinault¹, Xing-zhi Song¹, Yue Liu¹, Ye Yuan¹, Lynne Nazareth¹, Xiang Qin¹, Donna M. Musty¹, Marcel Margulies¹, George M. Weinstock¹, Richard A. Gibbs¹ & Jonathan M. Rothberg¹



Nature (2008)

A highly annotated whole-genome sequence of a Korean individual

Jong-Il Kim^{1,2,3,4*}, Young Seok Ju^{1,2*}, Hansoo Park^{1,2}, Sheehyun Kim¹, Seonwook Lee¹, Jae-Hyuk Yi¹, Joann Mudge⁵, Neil A. Miller⁶, Dongwan Hong⁷, Callum J. Bell⁸, Hye-Sun Kim¹, In-Soon Chung¹, Woo-Chung Lee¹, Ji-Sun Lee¹, Seung-Hyun Seo¹, Ji-Young Yoon¹, Hyun-Nyun Woo¹, Heerook Lee¹, Donghwan Suh^{1,2}, Seungbok Lee^{1,2}, Hyun-Jin Kim^{1,2}, Maryam Yavartanoo^{1,2}, Minhye Kwak^{1,2}, Ying Zheng^{1,2}, Mi Kyeong Lee¹, Hyunjun Park¹, Jeong Yeon Kim¹, Omer Gokcumen¹, Ryan E. Mills¹, Alexander Wait Zaranek¹, Joseph Thakurta¹, Xiaodi Wu¹, Ryan W. Kim¹, Jim J. Hurley¹, Shujun Luo¹, Gary P. Schroth¹, Thomas D. Wu¹, HyunBin Kim¹, Kap-Seok Yang¹, Woong-Yang Park^{1,2}, Hyungtae Kim¹, George M. Church¹, Charles Lee¹, Stephen F. Kingsmore¹ & Jeong-Sun Seo^{1,2,3,4,5}

Nature (2009)

1000 Genomes - Home

1000 Genomes

A Deep Catalog of Human Genetic Variation

Home About Partners Data Contact Wiki

1000 GENOMES PROJECT DATA RELEASE

SNP data downloads and genome browser representing four high coverage individuals

The first set of SNP calls representing the preliminary analysis of four genome sequences are now available to download through the EBI FTP site and the NCBI FTP site. The README file dealing with the FTP structure will help you find the data you are looking for.

The data can also be viewed directly through the 1000 Genomes browser at <http://browser.1000genomes.org>. Launch the browser and view a sample region here.

More information about the data release can be found in the data section of this web site.

Download the 1000 Genomes Browser Quick Start Guide

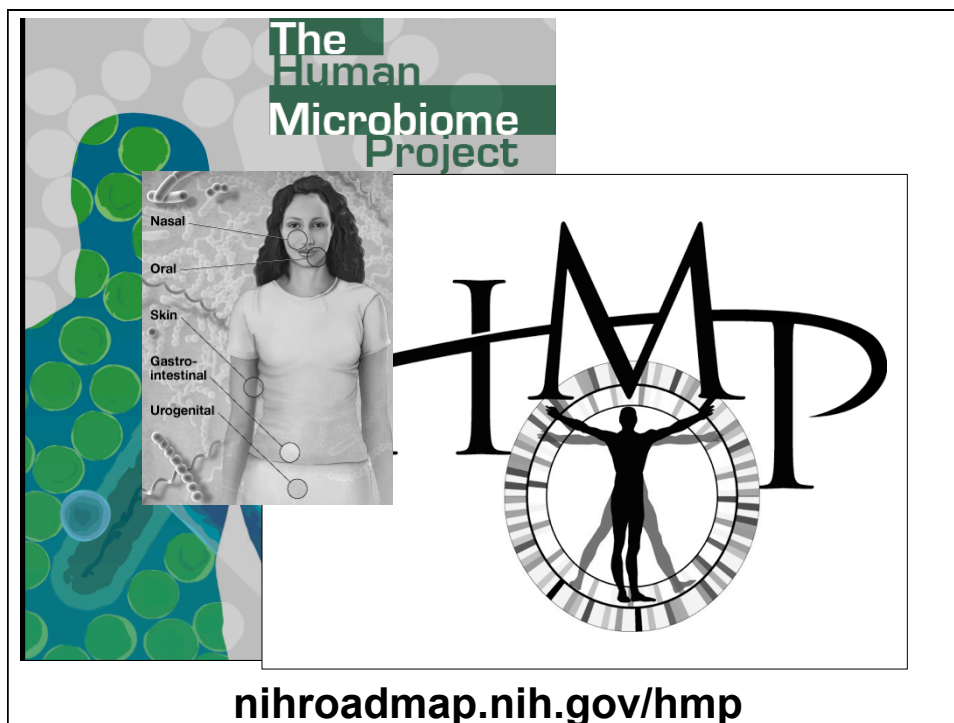
Quick start (pdf)

1000genomes.org



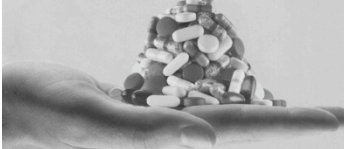
The screenshot shows the homepage of the The Cancer Genome Atlas (TCGA) website. At the top, it features the logos for the National Cancer Institute and the National Human Genome Research Institute. The main heading is "THE CANCER GENOME ATLAS" with a search bar and a "GO" button. Below the heading is a navigation menu with links for "About TCGA", "Program Components", "Policies", "Media Center", and "Launch Data Portal". The main content area includes a "Mission and Goal" section, a "TCGA Data Portal" section with links to "Access TCGA Data Portal" and "New* View the list of target genes and miRNAs selected by the TCGA network for sequencing", and a "News from the Pilot Project" section with a link to "NEW* TCGA Network Identifies More Than 6,000 Targets for Sequencing".

cancergenome.nih.gov



The graphic features the title "The Human Microbiome Project" in a stylized font. On the left, there is a silhouette of a human figure with various body sites labeled: "Nasal", "Oral", "Skin", "Gastro-intestinal", and "Urogenital". On the right, there is a large, stylized logo for "HMP" where the letters are interconnected, and a silhouette of a person stands within a circular element of the logo.

nihroadmap.nih.gov/hmp



All of these work.

Just not for everyone.

Perlegen may be able to help you sort out which medicine helps which patient.

Working with you, we can comprehensively analyze the DNA from thousands of patients taking your drug. Out of the millions of genetic variations between patients, we may be able to help you identify the ones that are associated with strong efficacy, poor efficacy, or side effects.

Perlegen's exceptional coverage of the genome and experienced team of analysts could help you get clinically relevant answers, not just data, in a matter of months.

We partner with the top pharmaceutical companies around the world. We also license late-stage drugs. If you have a drug that can benefit from our approach, please contact us.

Patients are waiting.

genetic@perlegen.com
Mountain View, California • 650-625-4500
Tokyo, Japan • 81 (0)3 3444-6080
www.perlegen.com

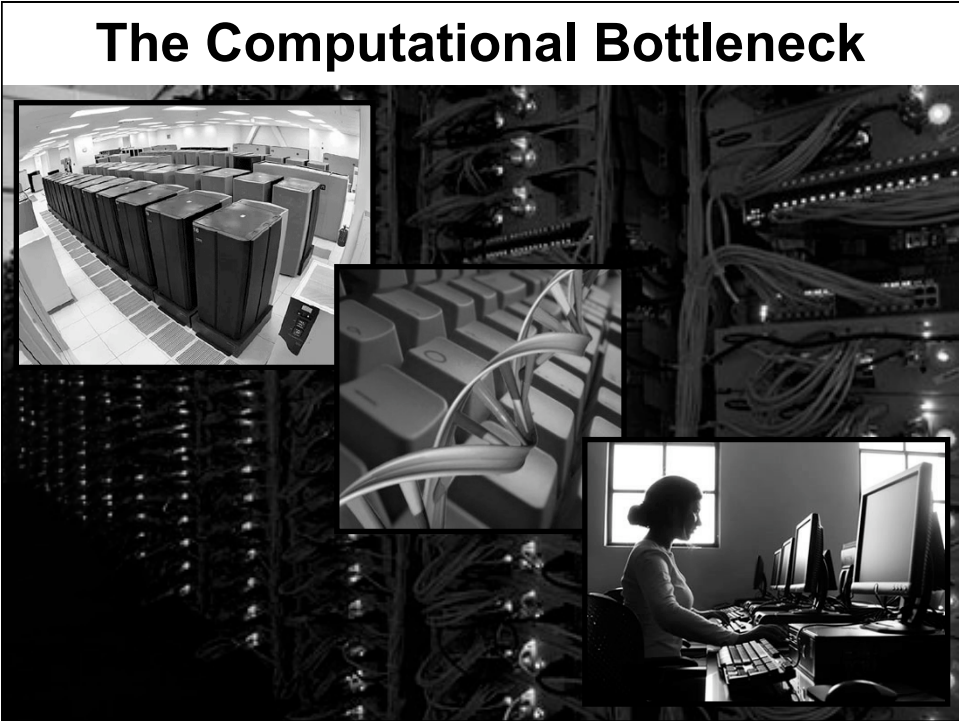
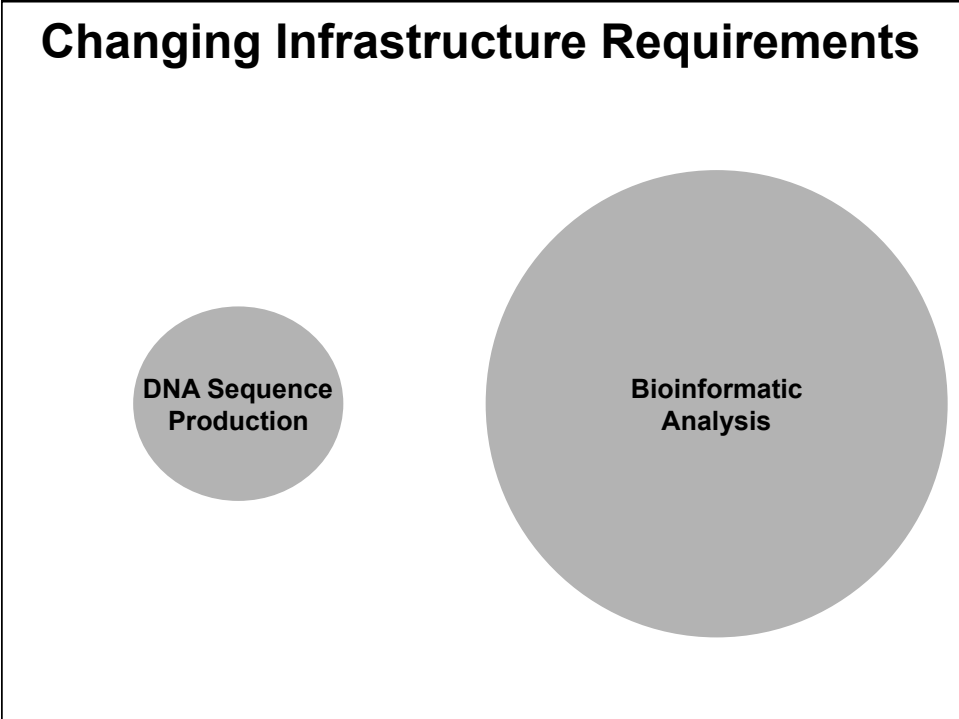
Targeting today's drugs.
Discovering tomorrow's.™

PERLEGEN
SCIENCE IS

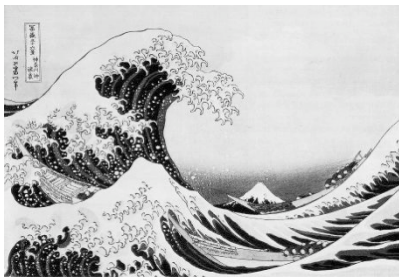
© 2009 Perlegen

Realities of New DNA Sequencing Technologies...





The Genomic Era: *circa* 2010



The Genomic Revolution Continues

The Top 10 Medical Advances of the Decade

By PEGGY PECK and LAUREN COX
 ABC News Medical Unit in Collaboration with MedPage Today
 Dec. 17, 2009

NEWS
 6 comments

Share this story with friends
 Digg submit Facebook Twitter Reddit StumbleUpon More

The first decade of the 21st century brought a number of discoveries, mistakes and medical advances that influenced medicine from the patient's bedside to the medicine cabinet.



Dr. John Sulston, Director of the Sanger Centre near Cambridge takes part in the Human Genome Project. (No New/Reuters)

In some cases, these advances changed deeply rooted beliefs in medicine. In others, they opened up possibilities beyond what doctors thought was possible years ago.

ABC News, in collaboration with MedPage Today, reached out to more than 800 specialists for their suggestions. More than 125 experts in various fields and specialties responded. Their suggestions were then sent to the American Association for the History of Medicine, which narrowed the pool down to an authoritative list of 10 medical advances this decade that have had the most impact.

1. Human Genome Discoveries Reach the Bedside

In 2000, scientists in California released a rough draft of the human genome to the public on the Internet. For the first time, the world could download and read the complete set of human genetic information and begin to discover what our roughly 23,000 genes do.

Mapping the human genome was a race involving time and money in the 1990s, with two competitors at the lead -- the government-funded Human Genome Project, which completed its task in 15 years using more than \$3 billion in taxpayer money, and a private company, Celera Genomics, which used \$100 million and took less than a decade.

Both groups announced drafts of the human genome at a June 26, 2000 press conference with then president Bill Clinton and former British Prime Minister Tony Blair.