

*Computational gene finding in
the human genome:
how many genes do we have?*

Steven Salzberg

Director, Center for Bioinformatics and Computational Biology

Horvitz Professor, Dept. of Computer Science

University of Maryland

<http://cbcb.umd.edu>

We've been trying for a long time to determine how many genes we have

letters to nature

Nature 201, 847 (22 February 1964); doi:10.1038/201847a0

A Preliminary Estimate of the Number of Human Genes

F. VOGEL

Institut für Anthropologie und Humangenetik, University of Heidelberg, Germany.

RECENT results of molecular genetics enable us to estimate the number of human genes, if certain assumptions are made. The following data are available: (1) The α -chain of human haemoglobin contains 141, the β -chain contains 146 amino-acids, corresponding to a molecular weight of about 17,000 each¹. Assuming a triplet code^{2,3} this means that the α - and β -chains are determined by 423 and 438 nucleotide pairs, respectively. According to 'Svedberg's law'⁴, many proteins consist of sub-units of the same order of magnitude (molecular weight of about 17,500). Hence, the assumption seems to be warranted that one average structural gene might have a length of about 450 nucleotide pairs. (2) The weight of one haploid human chromosome set in human spermatozoa is about 2.72×10^{-12} g. Granulocytes contain about 6.23×10^{-12} g; lymphocytes contain about 5.84×10^{-12} g (ref. 5). Extensive examinations have shown that the DNA content is constant in all resting cells of one species, which have the same number of chromosome sets, and depends on the degree of polyploidy^{5,6}. The assumption seems to be justified that most of the DNA works as genetic material, even if in some cells minor fractions with other functions might possibly be present⁷. In the following calculations the total amount of DNA in a haploid human chromosome set is estimated to be about 3×10^{-12} g. (3) Usually the genetic variants of human haemoglobins differ in one amino-acid substitution only^{1,8}. One structural gene can only produce one single type of genetically determined polypeptide chain. As much as we know, this applies for other genetically determined proteins as well. This means that the genetic information for these structural genes can only be present once. Any degree of polyteny for these loci in the germ cells is highly unlikely. As has been mentioned, however, the DNA content of diploid cells is about twice the content of (haploid) spermatozoa. We assume that the total genetic information is only present once.

1,000,000 genes? 100,000 genes?

Science 25 October 1996:

Vol. 274, no. 5287, pp. 540 - 546

DOI: 10.1126/science.274.5287.540

[< Prev](#) | [Table of Contents](#) | [Next >](#)

ARTICLES

A Gene Map of the Human Genome

G. D. Schuler, M. S. Boguski, E. A. Stewart, L. D. Stein, G. Gyapay, K. Rice, R. E. White, P. Rodriguez-Tomé, A. Aggarwal, E. Bajorek, S. Bentolila, B. B. Birren, A. Butler, A. B. Castle, N. Chiannikulchai, A. Chu, C. Clee, S. Cowles, P. J. R. Day, T. Dibling, N. Drouot, I. Dunham, S. Duprat, C. East, C. Edwards, J.-B. Fan, N. Fang, C. Fizames, C. Garrett, L. Green, D. Hadley, M. Harris, P. Harrison, S. Brady, A. Hicks, E. Holloway, L. Hui, S. Hussain, C. Louis-Dit-Sully, J. Ma, A. MacGilvery, C. Mader, A. Maratukulam, T. C. Matise, K. B. McKusick, J. Morissette, A. Mungall, D. Muselet, H. C. Nusbaum, D. C. Page, A. Peck, S. Perkins, M. Piercy, F. Qin, J. Quackenbush, S. Ranby, T. Reif, S. Rozen, C. Sanders, X. She, J. Silva, D. K. Slonim, C. Soderlund, W.-L. Sun, P. Tabar, T. Thangarajah, N. Vega-Czarny, D. Vollrath, S. Voyticky, T. Wilmer, X. Wu, M. D. Adams, C. Auffray, N. A. R. Walter, R. Brandon, A. Dehejia, P. N. Goodfellow, R. Houlgate, J. R. Hudson Jr., S. E. Ide, K. R. Iorio, W. Y. Lee, N. Seki, T. Nagase, K. Ishikawa, N. Nomura, C. Phillips, M. H. Polymeropoulos, M. Sandusky, K. Schmitt, R. Berry, K. Swanson, R. Torres, J. C. Venter, J. M. Sikela, J. S. Beckmann, J. Weissenbach, R. M. Myers, D. R. Cox, M. R. James, D. Bentley, P. Deloukas, E. S. Lander, T. J. Hudson

The human genome is thought to harbor 50,000 to 100,000 genes, of which about half have been sampled to date in the form of expressed sequence tags. An

Proceedings of the National Academy of Sciences, Vol 90, 11995-11999, Copyright © 1993 by National Academy of Sciences

ARTICLE

Number of CpG Islands and Genes in Human and Mouse

F Antequera and A Bird

Estimation of gene number in mammals is difficult due to the high proportion of noncoding DNA within the nucleus. In this study, we provide a direct measurement of the number of genes in human and mouse. We have taken advantage of the fact that many mammalian genes are associated with CpG islands whose distinctive properties allow their physical separation from bulk DNA. Our results suggest that there are ~45,000 CpG islands per haploid genome in humans and 37,000 in the mouse. Sequence comparison confirms that about 20% of the human CpG islands are absent from the homologous mouse genes. Analysis of a selection of genes suggests that both human and mouse are losing CpG islands over evolutionary time due to de novo methylation in the germ line followed by CpG loss through mutation. This process appears to be more rapid in rodents. Combining the number of CpG islands with the proportion of island-associated genes, we estimate that the total number of genes per haploid genome is ~80,000 in both organisms.

CORRESPONDENCE

Nature Genetics **8**, 114 (1994)
doi:10.1038/ng1094-114a

Predicting the total number of human genes

Francisco Antequera¹ & Adrian Bird¹

¹Institute of Cell and Molecular Biology, Darwin Building, University of Edinburgh, Scotland EH9 3JR, UK

NEWS AND VIEWS

Nature Genetics **7**, 345 - 346 (1994)
doi:10.1038/ng0794-345

How many genes in the human genome?

Chris Fields¹, Mark D. Adams¹, Owen White¹ & J. Craig Venter¹

¹The Institute for Genomic Research, 932 Clopper Road, Gaithersburg, Maryland 20878, USA

Estimates in *Nature Genetics* (2000)

 © 2000 Nature America Inc. • <http://genetics.nature.com>

letter

Estimate of human gene number provided by genome-wide analysis using *Tetraodon nigroviridis* DNA sequence

Hugues Roest Crolius, Olivier Jaillon, Alain Bernot, Corinne Dasilva, Laurence Bouneau, Cécile Fischer, Cécile Fizames, Patrick Wincker, Philippe Brottier, Francis Quétier, William Saurin & Jean Weissenbach

- 28,000 - 34,000
- Based on alignments to pufferfish (*Tetraodon nigroviridis*)

Estimates in *Nature Genetics* (2000)

letter

 © 2000 Nature America Inc. • <http://genetics.nature.com>

Analysis of expressed sequence tags indicates 35,000 human genes

Brent Ewing & Phil Green

- Based on expressed sequence tag (EST) alignments to human chromosome 22

Estimates in *Nature Genetics* (2000)

 © 2000 Nature America Inc. • <http://genetics.nature.com>

letter

Gene Index analysis of the human genome estimates approximately 120,000 genes

Feng Liang, Ingeborg Holt, Geo Pertea, Svetlana Karamycheva, Steven L. Salzberg & John Quackenbush

- Based on assemblies of ESTs
- Extrapolated to whole human genome

Estimates in *Nature Genetics* (2000)

 © 2000 Nature America Inc. • <http://genetics.nature.com>

corrections

Gene Index analysis estimates the human genome contains 120,000 genes

F. Liang *et al.*

Nature Genet. 25, 239–240 (2000).

- Correction to avoid over-counting immunoglobulin ESTs
- 81,000 genes, based on assemblies of ESTs
- 57,000 genes, extrapolating from Chr 21 and 22

The gene count guessing game

Science 19 May 2000:
Vol. 288, no. 5469, pp. 1146 - 1147
DOI: 10.1126/science.288.5469.1146

[< Prev](#) | [Table of Contents](#) | [Next >](#)

NEWS OF THE WEEK

HUMAN GENOME PROJECT: And the Gene Number Is ...?

Elizabeth Pennisi

COLD SPRING HARBOR, NEW YORK--Even though a draft sequence of the human genome is nearing completion, biologists still don't know how many genes it contains. Indeed, the range of estimates seems to be growing rather than shrinking. The question lies at the core of our understanding of genetic complexity. If genomes are the books of life, then genes are the words that tell the story of each organism. Biologists have long assumed that microorganisms are short stories and complex organisms such as humans, great tomes.



Place your bet. Uncertainty over the number of human genes has sparked a debate--and a betting pool.

Human genome paper I: *Nature* 409(15 Feb 2001), 860-921

articles

Initial sequencing and analysis of the human genome

International Human Genome Sequencing Consortium*

** A partial list of authors appears on the opposite page. Affiliations are listed at the end of the paper.*

- 30,000 - 40,000 genes
- Large degree of uncertainty about total
- Number of distinct transcripts and proteins even less certain

Human genome paper II: *Science* 291(16 Feb 2001), 1304-51

The Sequence of the Human Genome

J. Craig Venter,^{1*} Mark D. Adams,² Eugene W. Myers,¹ Peter W. Li,¹ Richard J. Mural,² Granger G. Sutton,³ Hamilton O. Smith,¹ Mark Yandell,¹ Cheryl A. Evans,¹ Robert A. Holt,² Jeanine D. Gocayne,¹ Peter Amanatides,¹ Richard M. Ballew,² Daniel H. Huson,¹ Jennifer Russo Wortman,¹ Qing Zhang,¹ Chinnappa D. Kodira,³ Xiangjun H. Zheng,¹ Lin Chen,¹ Marian Skupski,¹ Gangadharan Subramanian,¹ Paul D. Thomas,¹ Jinghui Zhang,¹ George L. Gabor Miklos,² Catherine Nelson,² Samuel Broder,¹ Andrew G. Clark,⁴ Joe Nadeau,² Victor A. McKusick,⁵ Norton Zinder,⁷ Arnold J. Levine,⁷ Richard J. Roberts,⁸ Mel Simon,⁸ Carolyn Slayman,¹² Michael Hunkapiller,¹¹ Randall Bolanos,² Arthur Delcher,² Ian Dew,¹ Daniel Fasulo,¹ Michael Flanigan,¹ Liliana Horea,¹ Aaron Halpern,¹ Sridhar Hannenhalli,¹ Saul Kravitz,¹ Samuel Levy,¹ Clark Mobarry,¹ Knut Reinert,¹ Karin Remington,¹ Jane Abu-Threideh,¹ Ellen Beasley,¹ Kendra Biddick,¹ Vivien Bonazzi,¹ Rhonda Brandon,¹ Michele Cargill,¹ Ishwar Chandramouliswaran,¹ Rosane Charlab,¹ Kabir Chaturvedi,¹ Zuoming Deng,¹ Valentina Di Francesco,² Patrick Dunn,² Karen Eilbeck,¹ Carlos Evangelista,¹ Andrei E. Gabrielian,¹¹ Weiniu Gan,² Wangmao Ge,¹ Fangcheng Gong,¹ Zhiping Gu,¹ Ping Guan,¹ Thomas J. Heiman,¹ Maureen E. Higgins,¹ Rui-Ru Ji,¹ Zhaoxi Ke,¹ Karen A. Ketchum,¹ Zhongwu Lai,¹ Yiding Lei,¹ Zhenya Li,¹ Jayin Li,¹ Yong Liang,¹ Xiaoying Lin,¹ Fu Lu,¹ Cannady V. Markovits,¹ Natalia Milshina,¹ Helen M. Moore,¹ Ashwinikumar K. Naik,¹ Vaibhav A. Narayan,¹ Beena Neelam,¹ Deborah Nusskern,¹ Douglas B. Rusch,¹ Steven Salzberg,¹² Wei Shao,¹ Bixiong Shue,¹ Jingtao Sun,¹ Zhen Yuan Wang,¹ Aihui Wang,¹ Xin Wang,¹ Jian Wang,¹ Ming-Hui Wei,¹ Ron Wildes,¹² Chunlin Xiao,¹ Chunhua Yan,¹ Alison Yao,¹ Jane Ye,¹ Ming Zhan,¹ Weiqing Zhang,¹ Hongyu Zhang,¹ Qi Zhao,¹ Liansheng Zheng,¹ Fei Zhong,¹ Wenyang Zhong,¹ Shiaoqing C. Zhu,¹ Shaying Zhao,¹² Dennis Gilbert,¹ Suzanna Baumhueter,¹ Gene Spier,¹ Christine Carter,¹ Anibal Cravchik,¹ Trevor Woodage,¹ Feroze Ali,¹ Huijin An,¹ Aderonke Awo,¹ Danita Baldwin,¹ Holly Baden,¹ Mary Barnstead,¹ Ian Barrow,¹ Karen Beeson,¹ Dana Busam,¹ Amy Carver,¹ Angela Center,¹ Ming Lai Cheng,¹ Liz Curry,¹ Steve Donaher,¹ Lionel Davenport,¹ Raymond Deslats,¹ Susanne Dietz,¹ Kristina Dodson,¹ Lisa Doup,¹ Steven Ferreira,¹ Neha Garg,¹ Andres Gluecksmann,¹ Brit Hart,¹ Jason Haynes,¹ Charles Haynes,¹ Cheryl Heiner,¹ Suzanne Hladun,¹ Damon Hostin,¹ Jarrett Houck,¹ Timothy Howland,¹ Chinyere Ibegwam,¹ Jeffery Johnson,¹ Francis Kalush,¹ Lesley Kline,¹ Shashi Koduru,¹ Amy Love,¹ Felecia Mann,¹ David May,¹ Steven McCawley,¹ Tina McIntosh,¹ Ivy McMullen,¹ Mee Moy,¹ Linda Moy,¹ Brian Murphy,¹ Keith Nelson,¹ Cynthia Pfannkoch,¹ Eric Prots,¹ Vinita Puri,¹ Hina Qureshi,¹ Matthew Reardon,¹ Robert Rodriguez,¹ Yu-Hui Rogers,¹ Deanna Romblad,¹ Bob Ruhfel,¹ Richard Scott,¹ Cynthia Sitter,¹ Michelle Smallwood,¹ Erin Stewart,¹ Renee Strong,¹ Ellen Suh,¹ Reginald Thomas,¹ Ni Ni Tint,¹ Sukyoo Tse,¹ Claire Vech,¹ Gary Wang,¹ Jeremy Wetter,¹ Sherita Williams,¹ Monica Williams,¹ Sandra Windsor,¹ Emily Winn-Den,¹ Keriolien Wolfe,¹ Jayshree Zaveri,¹ Karena Zaveri,¹ Josep F. Abril,¹⁴ Roderic Guigé,¹⁴ Michael J. Campbell,¹ Kimmen V. Sjolander,¹ Brian Karlak,¹ Anith Kejarwal,¹ Huaiyu Mi,¹ Betty Lazareva,¹ Thomas Hatton,¹ Apurva Narechania,¹ Karen Diemer,¹ Anushya Muruganujan,¹ Nan Guo,¹ Shinji Sato,¹ Vineet Bafna,¹ Sorin Istrail,¹ Ross Uppert,¹ Russell Schwartz,¹ Brian Walenz,¹ Shibo Yooseph,¹ David Allen,¹ Anand Basu,¹ James Baxendale,¹ Louis Blick,¹ Marcelo Caminha,¹ John Carnes-Stine,¹ Parris Caulk,¹ Yen-Hui Chiang,¹ My Coyne,¹ Carl Dahlke,¹ Anne Deslattes Mays,¹ Maria Dombroski,¹ Michael Donnelly,¹ Dale Ely,¹ Shiva Esparham,¹ Carl Foster,¹ Harold Gire,¹ Stephen Glanowski,¹ Kenneth Glasser,¹ Anna Glodok,¹ Mark Gorokhov,¹ Ken Graham,¹ Barry Gropman,¹ Michael Harris,¹ Jeremy Heil,¹ Scott Henderson,¹ Jeffrey Hoover,¹ Donald Jennings,¹ Catherine Jordan,¹ James Jordan,¹ John Kasha,¹ Leonid Kagan,¹ Cheryl Kraft,¹ Alexander Levitsky,¹ Mark Lewis,¹ Xiangjun Liu,¹ John Lopez,¹ Daniel Ma,¹ William Majoros,¹ Joe McDaniel,¹ Sean Murphy,¹ Matthew Newman,¹ Trung Nguyen,¹ Ngoc Nguyen,¹ Marc Nodell,¹ Sue Pan,¹ Jim Peck,¹ Marshall Peterson,¹ William Rowe,¹ Robert Sanders,¹ John Scott,¹ Michael Simpson,¹ Thomas Smith,¹ Arlan Sprague,¹ Timothy Stockwell,¹ Russell Turner,¹ Eli Venter,¹ Mei Wang,¹ Meiyuan Wen,¹ David Wu,¹ Mitchell Wu,¹ Ashley Xia,¹ Ali Zandieh,¹ Xiaohong Zhu,¹

- 26,588 genes
- 12,000 additional “likely” genes based on similarity to mouse or other evidence

Steven Salzberg,¹²

Human genome version 2.0: *Nature* 431 (21 October 2004)

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 result of this finishing process. The current genome sequence (Build 35) contains 2.85 billion nucleotides interrupted by only 341 gaps. It covers ~99% of the euchromatic genome and is accurate to an error rate of ~1 event per 100,000 bases. Many of the remaining euchromatic gaps are associated with segmental duplications 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. 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.

- 20,000 - 25,000 genes

How do we find genes?

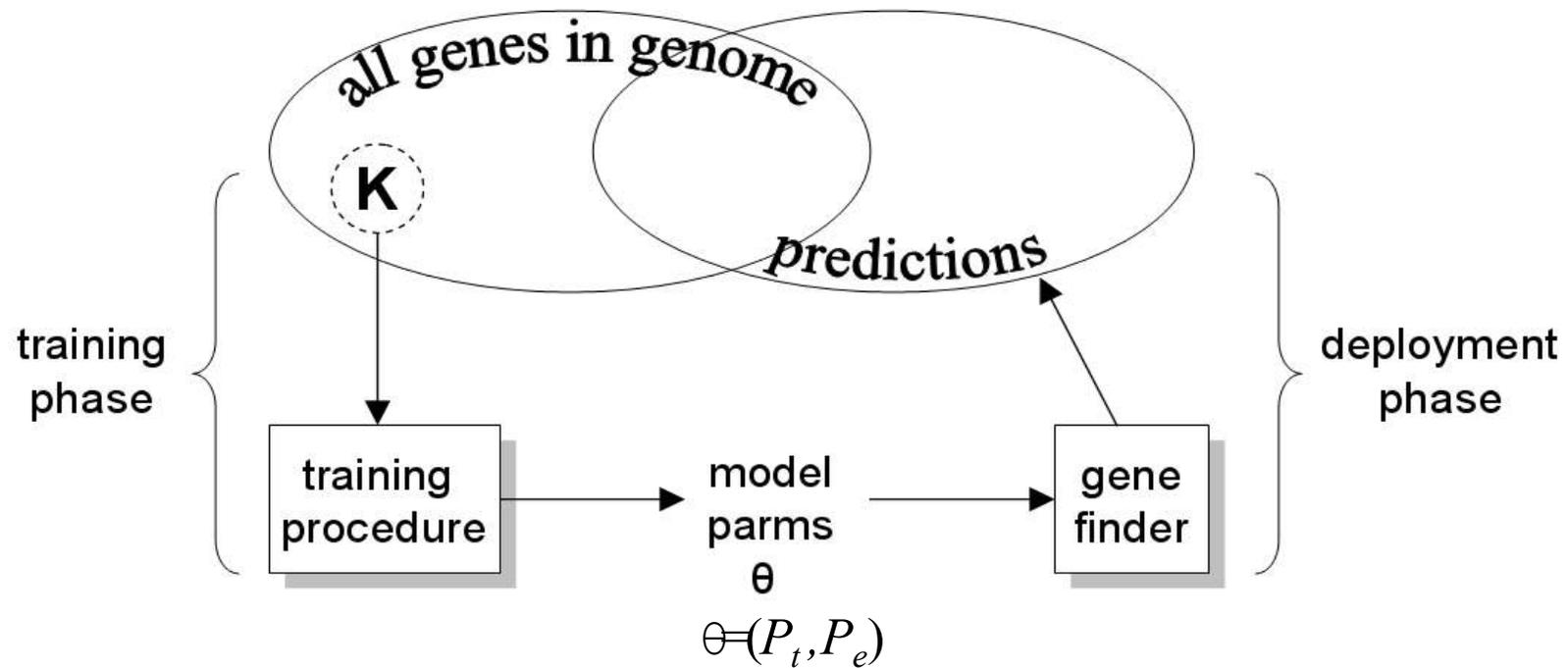
- *Ab initio* gene finding
- Expressed sequence tags (ESTs)
- Full-length cDNA sequencing
- Alignment of protein sequences to genomic DNA
- Combining all the evidence together







Training a Gene Finder



$$\theta_{\max} = \operatorname{argmax}_{\theta} \sum_{(S, \phi) \in K} P(\phi | S, \theta) = \operatorname{argmax}_{\theta} \frac{\sum_{(S, \phi) \in K} P_t(q_0, q_L) \prod_{i=1}^L P_e(x_i | q_i) P_t(q_i | q_{i-1})}{P(S | \theta)}$$

Search

SEARCH

GLIMMERHMM

Eukaryotic Gene-Finding System

University of Maryland » UMIACS » CBCB » Glimmer

OVERVIEW

GlimmerHMM is a new gene finder based on a Generalized Hidden Markov Model (GHMM). Although the gene finder conforms to the overall mathematical framework of a GHMM, additionally it incorporates splice site models adapted from the [GeneSplicer](#) program and a decision tree adapted from [GlimmerM](#). It also utilizes Interpolated Markov Models for the coding and noncoding models. Currently, GlimmerHMM's GHMM structure includes introns of each phase, intergenic regions, and four types of exons (initial, internal, final, and single). A basic user manual can be consulted [here](#).

SYSTEM REQUIREMENTS

GlimmerHMM is released as source code and was tested on Linux RedHat 6.x+, Sun Solaris, and Alpha OSF1, but should work on any Unix system.

ACCURACY

GlimmerHMM has been trained on several species including *Arabidopsis thaliana*, *Coccidioides* species, *Cryptococcus neoformans*, and *Brugia malayi*.
New: trainings for *C. elegans* and *Danio rerio* (zebrafish) are now available!

	<i>Nuc Sens</i>	<i>Nuc Spec</i>	<i>Nuc Accur</i>	<i>Exon Sens</i>	<i>Exon Spec</i>	<i>Exact Genes</i>	<i>Size of test set</i>
<i>D.rerio</i>	93%	78%	86%	77%	69%	24%	549 genes
<i>C.elegans</i>	96%	95%	96%	82%	81%	42%	1886 genes
<i>Arabidopsis</i>	97%	99%	98%	84%	89%	60%	809 genes
<i>Cryptococcus</i>	96%	99%	98%	86%	88%	53%	350 genes
<i>Coccidioides</i>	99%	99%	99%	84%	86%	60%	503 genes
<i>Brugia</i>	93%	98%	95%	78%	83%	25%	477 genes

GlimmerHMM has been recently trained on human. The table below presents its performance compared to Genscan on 963 human RefSeq genes selected randomly from all 24 chromosomes, non-overlapping with the training set. The test set contains 1000 bp of untranslated sequence on either side (5' or 3') of the coding portion of each gene.

	<i>Nuc Sens</i>	<i>Nuc Spec</i>	<i>Nuc Acc</i>	<i>Exon Sens</i>	<i>Exon Spec</i>	<i>Exon Acc</i>	<i>Exact Genes</i>
<i>GlimmerHMM</i>	86%	72%	79%	72%	62%	67%	17%
<i>Genscan</i>	86%	68%	77%	69%	60%	65%	13%

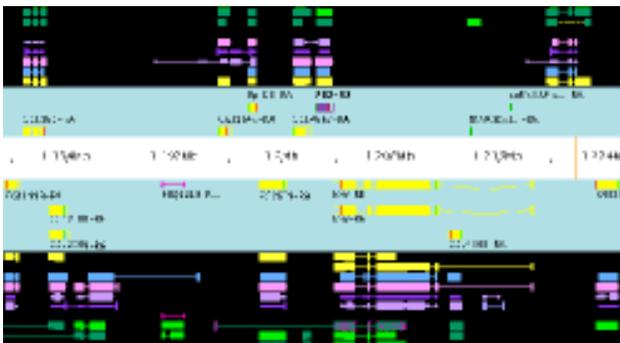
How do we find genes?

- *Ab initio* gene finding
- Expressed sequence tags (ESTs)
- Full-length cDNA sequencing
- Alignment of protein sequences to genomic DNA
- Combining all the evidence together

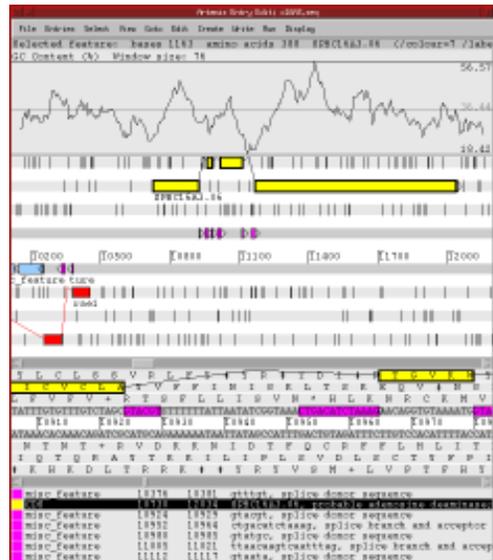
Putting it all together manually

- View *ab initio* predictions and sequence alignments within a genome editor:

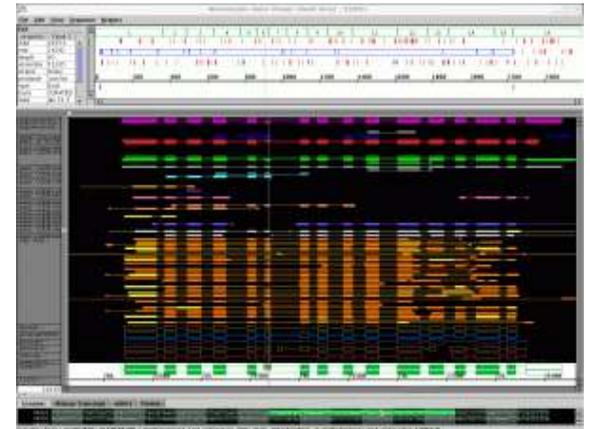
Apollo



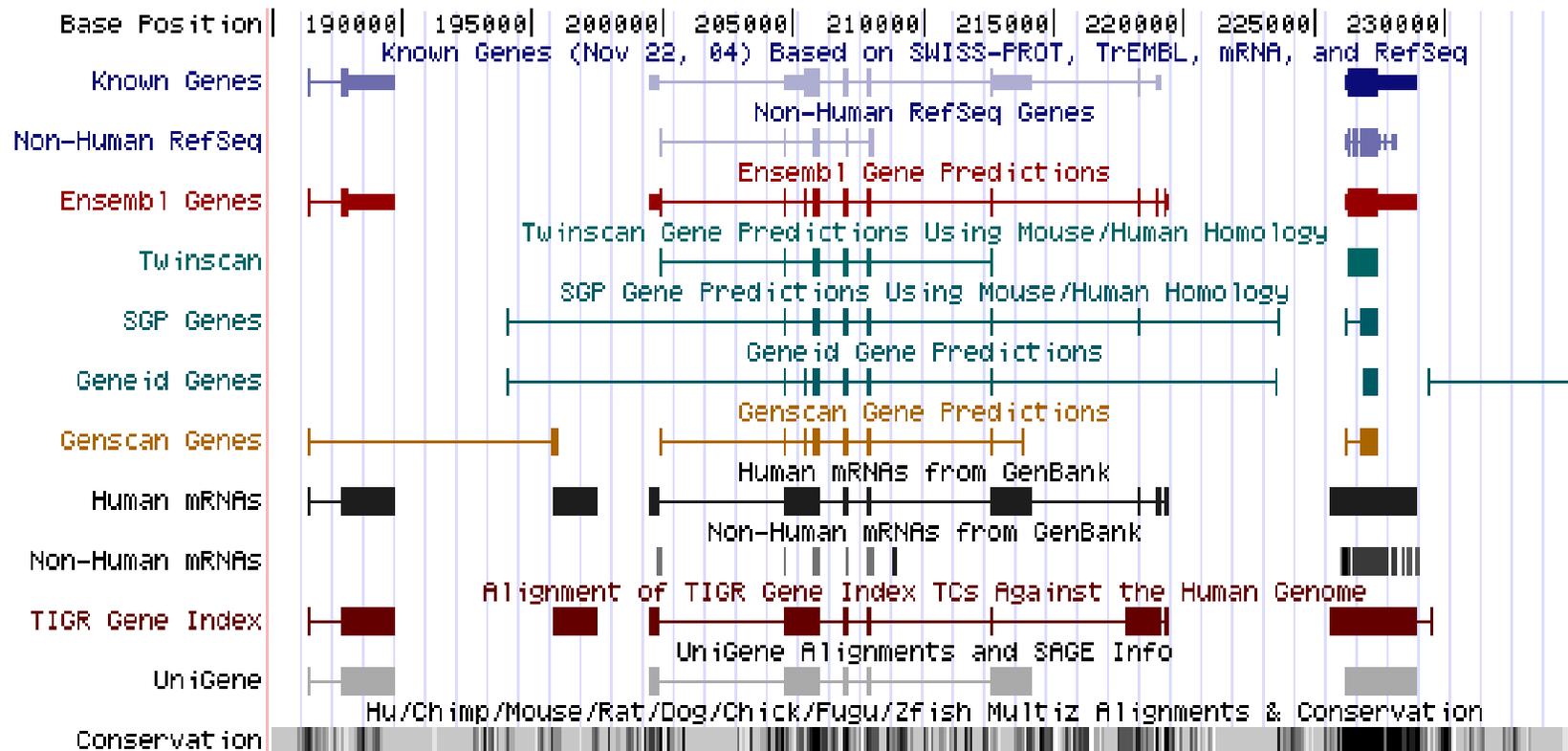
Artemis



Annotation Station



Putting it all together automatically with JIGSAW



Evaluating methods on 1% of the human genome: ENCODE

SPECIAL SECTION

GENES IN ACTION

VIEWPOINT

The ENCODE (ENCyclopedia Of DNA Elements) Project

The ENCODE Project Consortium*†

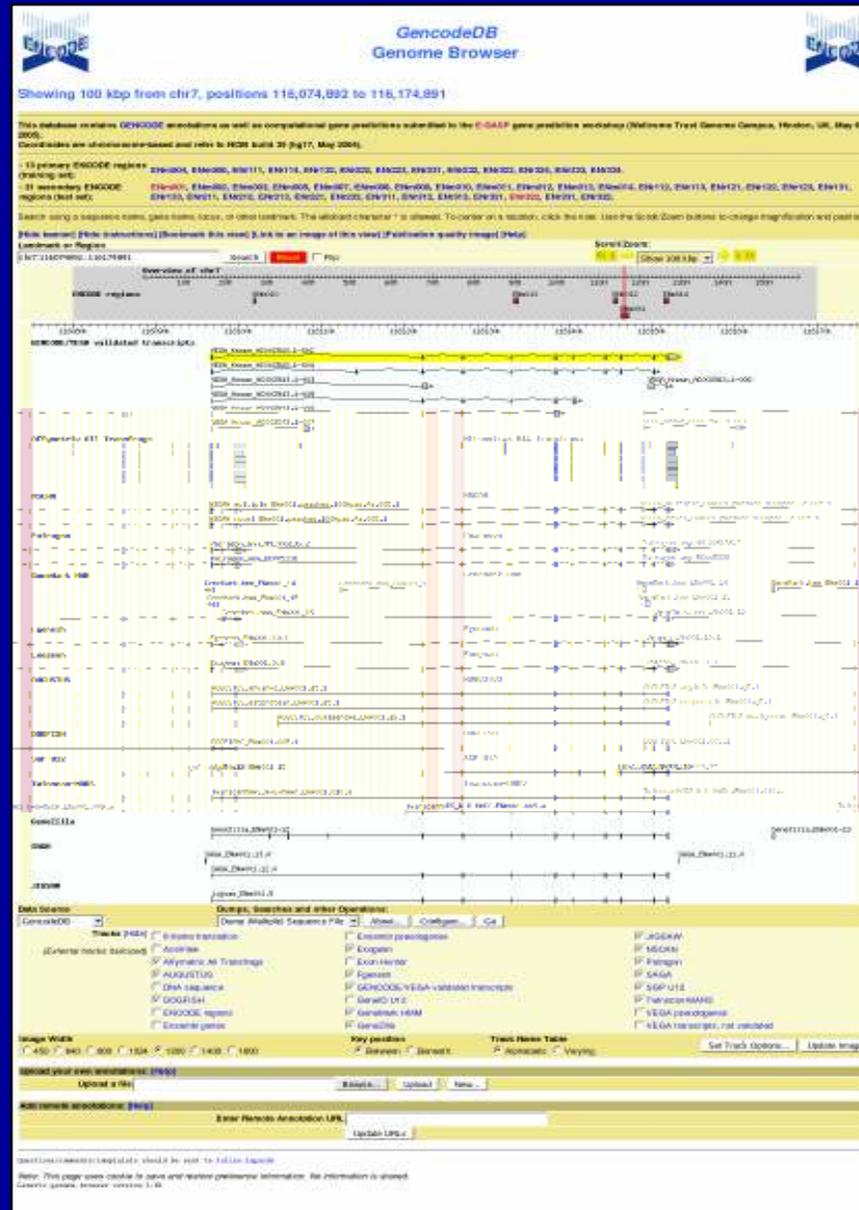
The ENCyclopedia Of DNA Elements (ENCODE) Project aims to identify all functional elements in the human genome sequence. The pilot phase of the Project is focused on a specified 30 megabases (~1%) of the human genome sequence and is organized as an international consortium of computational and laboratory-based scientists working to develop and apply high-throughput approaches for detecting all sequence elements that confer biological function. The results of this pilot phase will guide future efforts to analyze the entire human genome.

approaches, such as cDNA-cloning efforts (4, 5) and chip-based transcriptome analyses (6, 7), have revealed the existence of many transcribed sequences of unknown function. As a reflection of this complexity, about 5% of the human genome is evolutionarily conserved with respect to rodent genomic

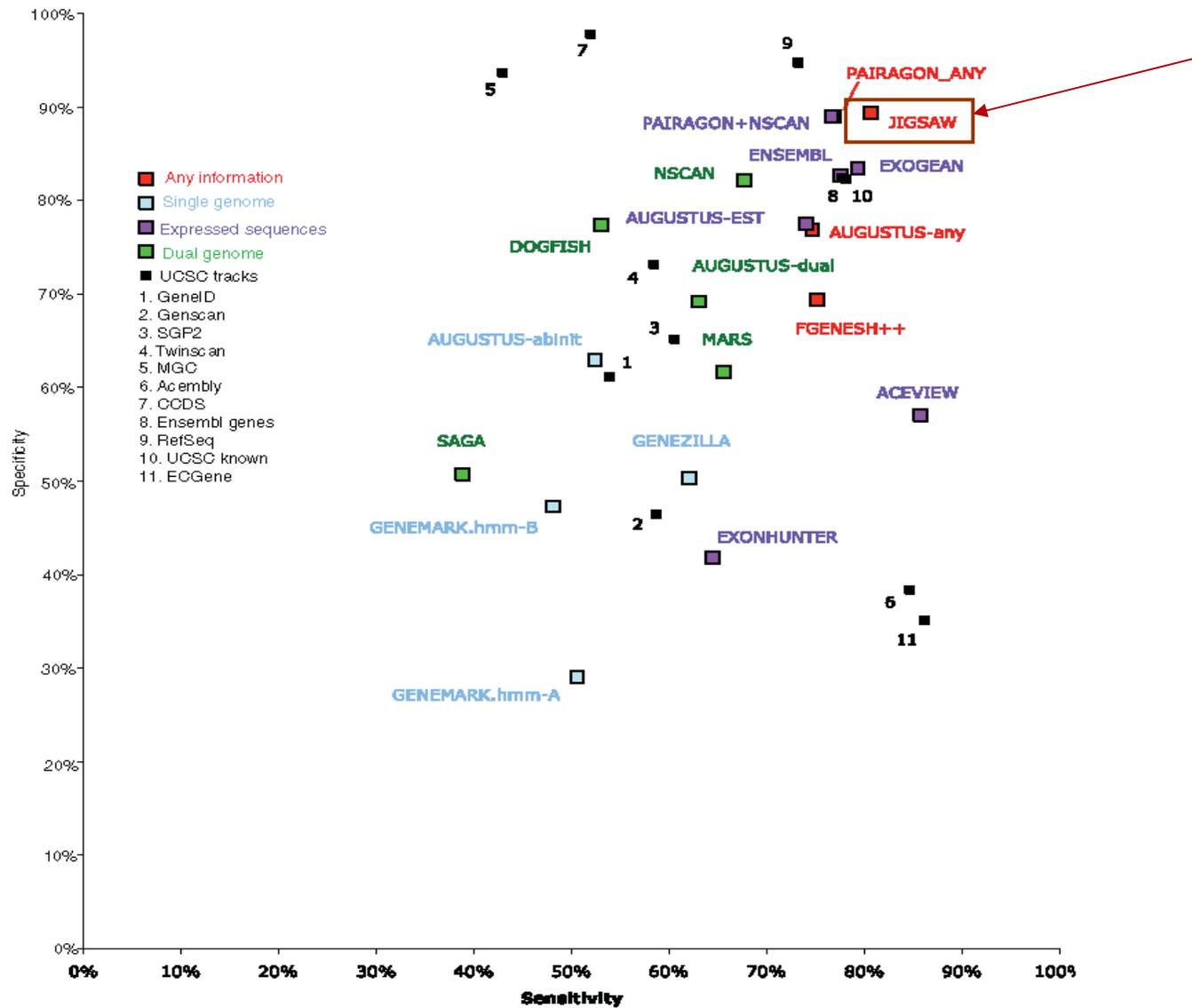
With the complete human genome sequence now in hand (1-3), we face the enormous challenge of interpreting it and learning how to use that information to understand the biology of human health and disease. The ENCyclopedia Of DNA Elements (ENCODE) Project is predicated on the belief that a comprehensive catalog of the structural and functional components encoded in the human genome sequence will be critical for understanding human biology well enough to address those fundamental aims of biomedical research. Such a complete catalog, or "parts list," would include protein-coding

Science 306 (2004), 636-640.

ENCODE Gene finding ASsessment Project (EGASP)



The GencodeDB Genome Browser. Fig. 2 from Guigo *et al.*, *Genome Biology* 2006, 7(Suppl 1):S2.



EGASP results: Exon prediction accuracy among 28 different methods

EGASP results: overall gene accuracy

<u>Genefinding Method</u>	<u>Sens</u>	<u>Spec</u>
AUGUSTUS-any	47.9	35.5
FGENESH++	69.9	42.0
JIGSAW	72.6	65.9
PAIRAGON-any	69.5	61.3
AUGUSTUS-abinit	24.3	17.2
GENEMARK.hmm-A	15.2	3.2
GENEMARK.hmm-B	16.8	7.9
GENEZILLA	19.5	8.8
ACEVIEW	63.5	48.6
AUGUSTUS-EST	47.6	37.0
ENSEMBL	71.6	67.3
EXOGEAN	63.1	80.8
EXONHUNTER	21.9	6.3
PAIRAGON+NSCAN_EST	69.5	61.7
AUGUSTUS-dual	26.0	18.6
DOGFISH	10.8	14.6
MARS	33.4	24.9
NSCAN	35.4	36.7
SAGA	4.3	3.4
GENSCAN	15.5	10.1
KNOWNgene	77.0	72.7
TWINSKAN	22.3	20.2

So, where are we now?

Ensembl genes
www.ensembl.org



- 21,774 “Known genes”
- 1,036 Novel genes
- 3,994 RNA genes
- 69,185 Genscan gene predictions

Let's not forget pseudogenes

- 27,130 pseudogenes
- www.pseudogene.org

Current NCBI gene counts

www.ncbi.nih.gov

- Entrez Gene:
 - 38,621 genes
 - But this includes pseudogenes
- RefSeq:
 - 28,961 genes
 - 31,784 transcripts

CCDS

- Home
- FTP
- Process
- Statistics

Collaborators

- EBI
- NCBI
- UCSC
- WTSI

Contact Us

- GenComp eMail

Genome Displays

- Ensembl
- Genome Browser
- Map Viewer
- VEGA

Related Resources

- Entrez Gene
- HomoloGene
- RefSeq
- UniGene

[Initial statistics for Homo sapiens for build 35.1](#)

[Current statistics for Homo sapiens for build 35.1](#)

[Initial statistics for Mus musculus for build 36.1](#)

[Current statistics for Mus musculus for build 36.1](#)

Initial statistics for *Homo sapiens* for build 35.1 ↑

as of March 2, 2005

CCDS Totals

Category	Count
CCDS IDs	14,795
Gene IDs	13,142
Sequence IDs	31,724

Feb. 26, 2007

Sequence IDs by Organization

NCBI RefSeq	15,496	19,360
EBI, WTSI Records	16,228	32,028

GeneID

Genes with >1 CCDS ID	1,205	1,708
-----------------------	-------	-------

Lesson:
Science isn't
decided by
voting

So, we still don't have a gene count

...and for many genes, we aren't yet sure of their exon-intron structure...

...and there are > 1000 other genomes already complete or under way...

...so we aren't giving up!

Thanks...



Mihaela Pertea



Jonathan Allen
(LLNL)



Bill Majoros
(Duke Univ.)



Art Delcher



Brian Haas
(Broad Institute)



United States
National Library of Medicine
National Institutes of Health

Evaluating Gene Predictions

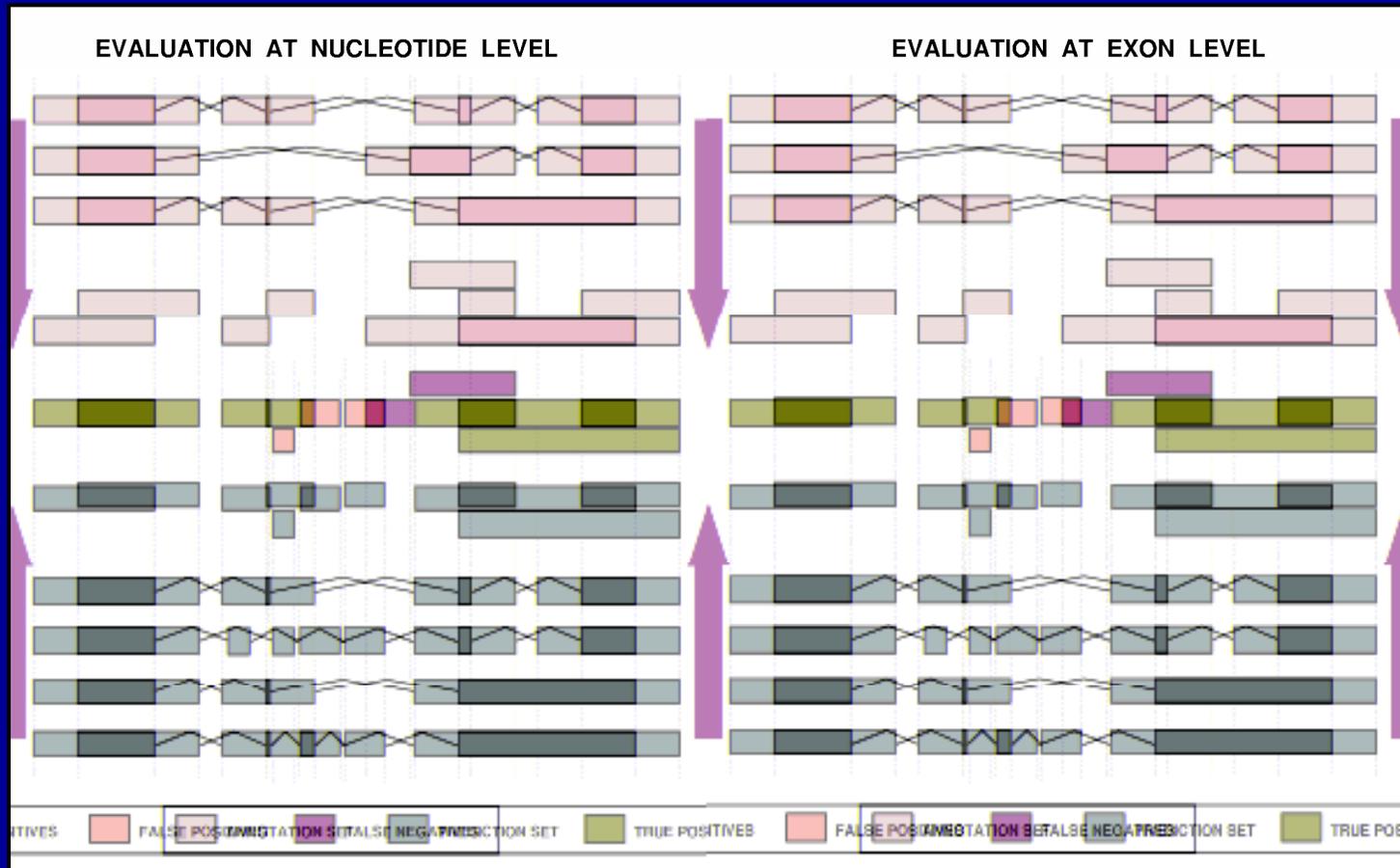


Fig. 3 from Guigo *et al.*, *Genome Biology* 2006, 7(Suppl 1):S2.