

TREBALL DE RECERCA
PER OPTAR AL
DIPLOMA D'ESTUDIS AVANÇATS
(DEA)

Josep Francesc ABRIL FERRANDO †

— 15 de març de 2002 —

Genome Informatics Research Lab
Grup de Recerca en Infomàtica Biomèdica
Institut Municipal d'Investigació Mèdica
Universitat Pompeu Fabra

† e-mail: **jabril@imim.es**

Índex

Articles Publicats	2
Visualitzant anotacions genòmiques amb el <code>gff2ps</code>	2
Avaluant l'anotació de la regió de 3Mbp <i>Adh</i> a <i>Drosophila melanogaster</i>	2
Anàlisi del genoma de <i>Drosophila melanogaster</i>	3
Anàlisi del genoma humà	3
Fiabilitat dels programes de predicció computacional de gens en seqüències genòmiques	4
Perspectives de Recerca	5
Predicció computacional de gens basada en l'homologia a nivell de seqüència	5
Anàlisi automatitzada del genoma humà.	5
Articles en preparació.	5
Appèndix.- Articles	7
<i>Bioinformatics</i> , 16(8):743 (2000)	7
<i>GenomeResearch</i> , 10(4):483 (2000)	8
<i>Drosophila</i> Genome Annotation Assessment Project	9
<i>Science</i> , 287(5461):2185 (2000)	10
Coding Content of the <i>Drosophila</i> Genome	11
<i>Science</i> , 291(5507):1304 (2001)	12
Annotation of the Celera Human Genome Assembly	13
<i>GenomeResearch</i> , 10(10):1631 (2000)	14

< Id: main.tex,v 1.3 2002/03/15 08:33:54 jabril Exp >

ARTICLES PUBLICATS

A continuació es presenta una llista dels articles publicats en els que entre els autors hi figura el doctorand. Primerament s'indica la pàgina d'aquest document on es poden veure incloses les dues primeres pàgines de cada article, ja que donat l'extensió d'alguns d'ells s'ha cregut convenient que es poden obtenir a través de la seva referència bibliogràfica; ara bé, el doctorand pot entregar una separata dels mateixos si fòra necessari. S'han inclòs també les figures principals corresponents a algun dels articles per destacar la feina en la que es va participar, aquestes figures estan escalades a la mida de la pàgina amb el que no es pot apreciar el detall de les mateixes, encara que es poden recuperar a través de la referència bibliogràfica corresponent.

Es pot accedir a més informació sobre les publicacions, incloent-hi els seus abstracts, així com els pòsters presentats a congressos científics, amb el corresponent *link* a la versió electrònica, des de la següent adreça d'internet:

<http://www1.imim.es/~jabril/index.html#PAPER>

★ Visualitzant anotacions genòmiques amb el gff2ps [pàgina 7]

Abril, J. F. and Guigó, R. (2000). gff2ps: visualizing genomic annotations. *Bioinformatics*, 16(8):743–744

Aquest programa genera, a partir d'un conjunt d'anotacions genòmiques en format GFF¹, una figura en format POSTSCRIPT on es resumeix de manera molt intuitiva la informació. Els avantatges de treballar amb aquest format de gràfics vectorials es poden resumir en els punts següents: és un format independent de dispositiu, al ser de tipus vectorial els gràfics es poden ampliar o reduir amb facilitat i permet treballar amb gran quantitat de dades de manera eficient. Malgrat que la seva publicació va ser a l'agost de l'any 2000, la primera versió operativa del programa ja era accessible a través d'internet des de principis de gener del mateix any, la qual cosa ens va permetre participar en la fase d'anotació del genoma de la mosca del vinagre (*Drosophila melanogaster*). En les següents tres publicacions es va participar directament en la representació gràfica de les anotacions disponibles per a les corresponents seqüències genòmiques. A més, el programa està sent utilitzat pel consorci públic que està seqüenciant i anotant el genoma de *Dictyostelium discoideum*², pels investigadors que estan desenvolupant el programa de predicción computacional de gens anomenat Phat³, etc...

Informació complementària:

➤ Adreça a internet del programa:

<http://www1.imim.es/software/gfftools/GFF2PS.html>

➤ Manual de l'usuari:

http://www1.imim.es/software/gfftools/gff2ps_docs/manual/MANUAL_GFF2PS_v0.96.ps.gz

★ Avaluant l'anotació de la regió de 3Mbp *Adh* a *Drosophila melanogaster* [pàgina 8]

Reese, M. G., Hartzell, G., Harris, N. L., Ohler, U., Abril, J. F., and Lewis, S. E. (2000). Genome annotation assessment in *Drosophila melanogaster*. *Genome Research*, 10:483–501

En aquest article es comparaven alguns dels mètodes computacionals de predicción de gens amb l'objectiu d'explorar la fiabilitat dels programes al treballar amb seqüències genòmiques d'organismes eucariotes. Els genomes eucariotes, a diferència dels procariotes, presenten la dificultat afegida de que la seva seqüència de DNA no codifica de manera contínua els gens i aquests a la

¹http://www.sanger.ac.uk/Software/formats/GFF/GFF_Spec.shtml

²<http://genome.imb-jena.de/dictyostelium/>

³<http://stat-www.berkeley.edu/users/scawley/Phat/README.html>

seva vegada tampoc estan codificats de manera contínua, però a més, la distribució dels gens al llarg de la seqüència tampoc és uniforme, ni tan sols les longituds dels introns. Fins la publicació d'aquest article no es disposava de seqüències llargues que continguessin múltiples gens i que haguéssin estat revisades amb cura per experts. Dotze grups, entre ells el nostre, van participar en el projecte d'avaluació d'anotacions genòmiques (GASP⁴) sobre un contig de 3Mbp del braç del cromosoma 2L de *D. melanogaster*, la regió on es troba el gen de la *Adh* i que conté uns 222 gens.

Informació complementària:

- La pàgina dedicada al pòster que es va presentar al ISMB'99:
<http://www1.imim.es/software/gfftools/GFF2PS-ADHposter.html>
- Versió electrònica del pòster, “*Drosophila* Genome Annotation Assessment Project”.
http://www.genome.org/content/vol10/issue4/images/data/483/DC1/GR10n4_poster.zip

★ Anàlisi del genoma de *Drosophila melanogaster* [pàgina 10]

Adams, M. D. *et al.* (including Josep F. Abril) (2000). The Genome Sequence of *Drosophila melanogaster*. *Science*, 287(5461):2185–2195

Aquest article conté l'anàlisi de la primera versió ensamblada del genoma de *D. melanogaster*, que van elaborar conjuntament el consorci públic de seqüenciació del genoma de la mosca, encapçalat pel *Berkeley Drosophila Genome Project* de la universitat de Berkeley a Califòrnia (BDGP⁵), i la empresa nordamericana Celera Genomics, dirigida pel Dr. Craig Venter. Va ser la prova de que el mètode de seqüenciació per fragments a l'atzar (*shotgun sequencing*) era aplicable a genomes d'organismes eucariotes complexos, alhora que era el genoma més gran seqüenciat fins a la data, 120Mbp de seqüència repartits en 4 cromosomes. L'anàlisi computacional d'aquest genoma va revelar l'existència d'uns 13000 gens, que es van representar en la figura central d'aquest article (veieu la figura 2, a la pàgina 11 d'aquest informe).

Informació complementària:

- Versió electrònica de la figura 4, “Coding content of the fly genome”.
<http://www.sciencemag.org/feature/data/genomes/2000/drosophila.shl>

★ Anàlisi del genoma humà [pàgina 12]

Venter, J. C. *et al.* (including J. F. Abril and R. Guigó) (2001). The Sequence of the Human Genome. *Science*, 291(5507):1304–1351

El febrer del 2001 es van publicar de manera simultània els dos esborranyos de la seqüència del genoma humà, l'un va aparèixer a la revista *Nature* i estava signat pel consorci públic, mentre que l'altre ho feia a la revista *Science* i venia signat per la empresa Celera Genomics. Donat que ja havíem col.laborat amb Celera en l'etapa d'anotació del genoma de *D. melanogaster*, ens van oferir l'oportunitat de poder treballar amb les dades del genoma humà de que disposaven. El genoma humà, de 3Gbp, era 30 vegades més gran que el de la mosca del vinagre, es treballava amb uns 35000 gens i tota la resta d'informació adicional (contingut en G+C, regions repetitives, etc...). Per complicar encara més les coses, la seqüència del genoma humà tan sols codifica en un 1%, la distribució dels gens al llarg de la seqüència dels cromosomes és molt variable així com la seva densitat (per exemple, el cromosoma 19 és un dels que té més gens per kilobase). Això ens va obligar a treballar amb dues escales, una per representar les regions on els gens mapaven sobre la seqüència dels cromosomes, i un'altra, a mena d'expansions que es projectaven a partir d'aquelles regions, per poder visualitzar les estructures exòniques d'aquests gens. Alhora, es va haver de paral·lelitzar l'execució del programa `gff2ps` (corrent el programa amb els mateixos

⁴<http://www.fruitfly.org/GASP1/>

⁵<http://www.fruitfly.org/>

paràmetres de configuració independentment per cada cromosoma i fussionant els resultats en una única figura) per poder obtenir esborranys de la figura en un temps raonable.

Informació complementària:

- Versió electrònica de la figura 1, “Annotation of the Celera Human Genome Assembly”.

<http://www.sciencemag.org/cgi/content/full/291/5507/1304/DC2>

★ **Fiabilitat dels programes de predicció computacional de gens en seqüències genòmiques [pàgina 14]**

Guigó, R., Agarwal, P., Abril, J. F., Burset, M., and Fickett, J. W. (2000). Gene prediction accuracy in large DNA sequences. *Genome Research*, 10:1631–1642

Donada la mancança de conjunts de dades per seqüències a escala genòmica, era necessari poder avaluar la fiabilitat i la precisió dels programes de predicció computacional de gens quan aquests s'enfrontessin amb aquestes dades. Aquest punt és força important en quant que el primer resultat que en resulta de seqüenciar el genoma humà és el catàleg de gens codificats en la seva seqüència. Per que aquest conjunt de gens tingui una utilitat a nivell biomèdic ha de ser el més acurat i complet que es pugui. Es va treballar amb dos conjunts de dades:

a) un extracte de la base de dades de seqüències EMBL, amb 178 seqüències genòmiques que contenen sencer un únic gen (i del qual es coneixen les coordenades del seu mRNA i dels exons que el composaven). La seva longitud mitjana era de 7Kbp.

b) sobre un fons de seqüència generat a l'atzar per simular els espais intergènics es van combinar conjunts de les seqüències anteriors. La longitud mitjana d'aquests pseudo-BACs era d'uns 180Kbp.

Els resultats suggerien que les eines de predicció computacional de gens a escala genòmica encara no són capaces de recuperar l'estructura exònica exacta de cadascun dels gens codificats en el genoma humà.

Informació complementària:

- El conjunt de seqüències de prova està disponible a:

<http://www1.imim.es/databases/gpeval2000/index.html>

PERSPECTIVES DE RECERCA

★ Predicció computacional de gens basada en l'homologia a nivell de seqüència

L'anàlisi comparativa de genomes permetrà una millora de les prediccions computacionals dels gens codificats en els mateixos. Emprant la informació que ens aporta l'homologia a nivell de seqüència entre espècies filogenèticament properes es poden abordar dos problemes que avui en dia encara no s'han resolt i que estan força interrelacionats: per una banda, completar el catàleg de gens per a una espècie en concret —en el nostre cas aplicant les tècniques per determinar les regions codificant del genoma humà—, i per l'altra, poder acostar-nos a la determinació del patró de processat alternatiu dels exons que conformen l'estructura d'un gen (el que s'anomena *alternative splicing*). La nostra aproximació es basa en utilitzar les dades de seqüència del genoma de ratolí (*Mus musculus*), tant dels contigs provenints del procés de seqüenciació per *shotgun* com de les seqüències genòmiques ensamblades dels seus cromosomes —colaborant estretament amb el consorci públic del genoma del ratolí—, per reanotar el genoma humà.

★ Anàlisi automatitzada del genoma humà.

Per poder fer front als requeriments que sorgeixen de l'apartat anterior, ens cal desenvolupar una sèrie de programes que s'executin de manera automàtica cada cop que vagin sortint noves versions tant a nivell de dades de seqüència com de les anotacions disponibles. Això ens permetrà adaptar-nos a uns conjunts de dades que van canviant al llarg del temps i, a la vegada, incrementar la fiabilitat dels nostres resultats al minimitzar les errades humanes que poden apareixer en escalar les aplicacions a un problema tan complexe com la reanotació d'un genoma.

★ Articles en preparació.

La feina que s'està duent a terme actualment es centra en la implementació de diferents programes per fer manipular i analitzar les dades d'origen genòmic, en concret:

- **gff2aplot**: visualitzant anotacions sobre alineaments de parells de seqüències (*pairwise alignments*). A partir d'un o més fitxers en format GFF⁶, que contenen informació tant de l'homologia entre parells de seqüències com de les anotacions per cadascuna de les seqüències, es genera un fitxer en POSTSCRIPT. Els principals avantatges respecte altres aplicacions actualment disponibles per visualitzar alineaments, són: que pot ser fàcilment incorporat en els programes que emprem per automatitzar tasques —donat que treballem amb el format GFF per intercanviar dades entre diferents tasques—, que pot incorporar molta informació sobre les anotacions de les seqüències analitzades i per la gran qualitat que ens proporcionen els gràfics vectorials amb format POSTSCRIPT.
- **SGP-2**: de *Syntenic Gene Prediction tool*, aquesta eina combina la informació de l'homologia entre seqüències de dues espècies que s'obté amb el programa TBLASTX —que forma part de la aplicació BLAST desenvolupada a Washington University⁷—, amb els algoritmes de predicció computacional de gens que implementa el programa geneid^{8,9}. La seva potència es basa en millorar la puntuació dels exons predictius que solapen regions conservades (els anomenats HSPs o *High-scoring Segment Pairs*), tot i projectant la puntuació d'aquests per aprofitar la seva informació, afegir aquestes puntuacions a les calculades pel geneid per cada exò abans d'agrupar-los dintre d'estructures gèniques.

⁶Veieu la nota a peu de pàgina número 1 de la pàgina 2

⁷Altschul, S. F., Gish, W., Miller, W., Myers, E. W., and Lipman, D. (1990). Basic local alignment search tool. *Journal of Molecular Biology*, 215:403–410

⁸Guigó, R., Knudsen, S., Drake, N., and Smith, T. F. (1992). Prediction of gene structure. *Journal of Molecular Biology*, 226:141–157

⁹Parra, G., Blanco, E., and Guigó, R. (2000). Geneid in *Drosophila*. *Genome Research*, 10:511–515

Un cop tinguem versions plenament operatives de tots dos programes, s'enviarà per la seva publicació un article per cadascun d'ells. Tots dos articles estan en aquests moments en fase de preparació; en el cas del SGP-2 l'article corresponent detallarà els algoritmes i els paràmetres amb els que haurem treballat sobre els conjunts de seqüències de prova.

APPÈNDIX.- ARTICLES

BIOINFORMATICS APPLICATIONS NOTE Vol. 16 no. 8/2000
Pages 743-744

gff2ps: visualizing genomic annotations

Josep F. Abril* and **Rodríguez Guijó**

Grup de Recerca en Informàtica Mèdica, Institut Municipal d'Investigació Mèdica (IMIM), Universitat Pompeu Fabra (UPF), C/ Dr. Aiguader, 80, 08003—Barcelona, Spain

Received on December 15, 1999; revised on February 18, 2000; accepted on February 24, 2000

Summary: *gff2ps* is a program for visualizing annotated features on a genomic sequence in GFF format as input, and produces a visual output in PostScript. While it can be used in a very simple way, it also allows for a great degree of customization through a number of options and/or customization files.

Availability: *gff2ps* is freely available at <http://www1.inim.es/~jabril/GFFTOOLS/GFF2PS.html>

Contact: *jabril@inim.es*

Supplementary information: <http://www1.inim.es/~jabril/GFFTOOLS/GFF2PS.html>

Abstract: *gff2ps* plots the features from different sources specified on a GFF file in a number of parallel rows (the so-called tracks here) along the length of the output pages (see Figure 1). Actually these are 'virtual' pages (the so-called blocks here) allowing for several blocks to be included in a single physical page, or for splitting a single block in a number of physical pages. Features can be plotted in a variety of colors and shapes and those grouped together can be visually linked in a number of ways.

gff2ps allows for a substantial amount of customization through command line options, and configuration files. However, meaningful output in most cases, meaningful output can be obtained without the need of any customization, by simply calling *gff2ps* with the input GFF file. *gff2ps* assumes, by default, that the GFF file itself carries enough formatting information. The examples in the figure show the versatility of *gff2ps*. Additional examples can be found at the *gff2ps* web page, as well as a detailed User Manual.

One of the main advantages of *gff2ps* is its ability to manage many physical page formats, including user-defined ones. This allows, for instance, the generation of poster size genomic maps. As an example, we used *gff2ps* to display at the ISMB'99 meeting, the predictions submitted to the Genome Annotation Assessment Project (GASP1) (<http://www.fruitfly.org/GASP1/>). The GASP1 plot was generated on three B0 size posters from a GFF file of over 50 000 feature records. The program has also been used to obtain the poster figures of recent relevant papers in genomic research (Adams *et al.*, 2000; Reecce, 2000).

Acknowledgments

We thank Moisés Burreci and Genís Parra (IMIM) for their useful comments. Richard Bruskiewich (Sanger Center) for his helpful hints on the GFF format, also Elena Casacuberta and Alfonso Monferr (CSIC) for motivating us to develop this tool. This work is supported by a grant from Plan Nacional de I+D, BIO98-0443-C02-01, and from a Fellowship to J.A. from the Instituto de Salud Carlos III, 99/9345.

The page description language PostScript is recognized as the current *de facto* industry standard for high-quality printing. PostScript provides both a printer-independent and a computer-system-independent means to describe integrated text and graphics, which can be put out on a variety of printers, plotters and workstation screens. The generation of PostScript output is very common in sequence analysis tools. Notably, we can cite the RSVP package by Searls (1993),

Fig. 1. Different views of the same input GFF file, using *gff2ps* with different configuration files and command-line options. The top two pages on the left were obtained using the default configuration (specifying only the number of pages, and output media size). By default, *gff2ps* makes a number of assumptions. Among others: (i) Features grouped from the GFF input file (ungrouped features are treated as a single element group) within the same source are plotted in the minimum number of tracks, guaranteeing that different groups do not overlap; (ii) The plot is fitted into a single block (assuring the length of the sequence to be the end of the most downstream feature), and the block is printed into a single physical page; (iii) Features for which the frame is specified are plotted using a two-color code scheme. The upstream half of the graphical element representing the frame of the feature and the downstream half the complement modulus 3 of its remainder. This is useful to check frame consistency between adjacent features (for instance, predicted exons). Two adjacent features are frame compatible when the color of the downstream half of the upstream feature matches the color of the upstream half of the downstream feature. (iv) If a score is provided for a feature, the feature is plotted with a height proportional to its score. (v) Obviously, all these default options can be overridden by the user. Notes: The real size color plots, the input GFF files, the configuration files and command line options used in each case, as well as additional examples can be found at: <http://www1.inim.es/~jabril/GFFTOOLS/GFF2PS-Shishplots.html>.

References

Adams, M.D. and Abril, J.F. *et al.* (2000) The genome sequence of *Drosophila melanogaster*. *Science*, **287**(5461), 2185-2195.

Searls, D.B. (1993) Doing sequence analysis with your printer. *Comput. Appl. Bioc.*, **9**(4), 421-426.

Lewis, S.E. (2000) Genome annotation assessment in *Drosophila melanogaster*. *Genome Res.*, **10**, 483-501.

*To whom correspondence should be addressed.

© Oxford University Press 2000

Genome Informatics Research Lab

Abril, JF

15 de març de 2002

Letter

Genome Annotation Assessment in *Drosophila melanogaster*

Martin G. Reese,^{1,4} George Hartzell,¹ Nomi L. Harris,¹ Uwe Ohler,^{1,2} Josep F. Abril,³ and Suzanna E. Lewis⁵

¹Berkeley Drosophila Genome Project, Department of Molecular and Cell Biology, University of California, Berkeley, California 94720-3200 USA; ²Chair for Pattern Recognition, University of Erlangen-Nuremberg, D-91058 Erlangen, Germany; ³Institut Municipal d'Investigació Médica—Universitat Pompeu Fabra, Departament of Medical Informatics (IMIM)—IUPF, 08003 Barcelona, Spain

Computational methods for automated genome annotation are critical to our community's ability to make full use of the large volume of genomic sequence being generated and released. To explore the accuracy of these automated feature prediction tools in the genomes of higher organisms, we evaluated their performance on a large, well-characterized sequence contig from the *Adh* region of *Drosophila melanogaster*. This experiment, known as the Genome Annotation Assessment project (GASP), was launched in May 1999. Twelve groups, applying state-of-the-art tools, contributed predictions for features, including gene structure, protein homologies, promoter sites, and repeat elements. We evaluated these predictions using two standards, one based on previously undisclosed high-quality full-length cDNA sequences and a second based on the set of annotations generated as part of an in-depth study of the region by a group of *Drosophila* experts. Although these standard sets only approximate the unknown distribution of features in this region, we believe that when taken in context the results of an evaluation based on them are meaningful. The results were presented as a tutorial at the conference on Intelligent Systems in Molecular Biology (ISMB'99) in August 1999. Over 95% of the coding nucleotides in the region were correctly identified by the majority of the gene finders, and the correct intron/exon structures were predicted for >90% of the genes. Homology-based annotation techniques recognized and associated functions with almost half of the genes in the region; the remainder were only identified by the ab initio techniques. This experiment presents the first assessment of promoter prediction techniques for a significant number of genes in a large contiguous region. We discovered that the promoter predictors' high false-positive rates make their predictions difficult to use. Integrating gene finding and DNA/EST alignments with promoter predictions decreases the number of false-positive classifications but recovers less than one-third of the promoters in the region. We believe that by establishing standards for evaluating genomic annotations and by assessing the performance of existing automated genome annotation tools, this experiment establishes a baseline that contributes to the value of ongoing large-scale annotation projects and should guide further research in genome informatics.

Genome annotation is a rapidly evolving field in genomics made possible by the large-scale generation of genomic sequences and driven predominantly by computational tools. The goal of the annotation process is to assign as much information as possible to the raw sequence of complete genomes with an emphasis on the location and structure of the genes. This can be accomplished by ab initio gene finding, by identifying homologues to known genes from other organisms, by the alignment of full-length or partial mRNA sequences to the genomic DNA, or through combinations of such methods. Related techniques can also be used to identify other features, such as the location of regulatory elements or repetitive sequence elements. The ultimate goal of genome annotation, the func-

tional classification of all the identified genes, currently depends on discovering homologies to genes with known functions.

We are interested in an objective assessment of the state of the art in automated tools and techniques for annotating complete genomes. The Genome Annotation Assessment Project (GASP) was organized to formulate guidelines and accuracy standards for evaluating computational tools and to encourage the development of new models and the improvement of existing approaches through a careful assessment and comparison of the predictions made by current state-of-the-art programs.

The GASP experiment, the first of its kind, was similar in many ways to the CASP (Critical Assessment of Techniques for Protein Structure Prediction) contests for protein structure prediction (Dunbrack et al. 1997;

Reese et al.

Levitt 1997; Moult et al. 1997, 1999; Sippl et al. 1999; Zemla et al. 1999), described at <http://predictioncenter.llnl.gov>. However, unlike the CASP contest, GASP was promoted as a collaboration to evaluate various techniques for genome annotation.

The GASP experiment consisted of the following stages: (1) Training data for the *Adh* region, including 2.9 Mb of *Drosophila melanogaster* genomic sequence, was collected by the organizers and provided to the participants; (2) a set of standards was developed to evaluate submissions while the participating groups produced and submitted their annotations for the region; and (3) the participating groups predictions were compared with the standards, a team of independent assessors evaluated the results of the comparison, and the results were presented as a tutorial at ISMB'99 (Reese et al. 1999).

Participants were given the finished sequence for the *Adh* region and some related training data, but they did not have access to the full-length cDNA sequences that were sequenced for the paper by Ashburner et al. (1999) that describes the *Adh* region in depth. The experiment was widely announced and open to any participants. Submitters were allowed to use any available technologies and were encouraged to disclose their methods. Because we were fortunate to attract a large group of participants who provided a wide variety of annotations, we believe that our evaluation addresses the state of art in genome annotation.

Twelve groups participated in GASP, submitting annotations in one or more of six categories: ab initio gene finding, promoter recognition, EST/cDNA alignment,

protein similarity, repetitive sequence identification, and gene function. Table 1 lists each participating group, the names of the programs or systems it used, and which of the six classes of annotations it submitted (see enclosed poster in this issue for a graphic overview of all the groups' results). Additional papers in this issue are written by the participants themselves and describe their methods and results in detail.

It should be noted that the lack of a standard that is absolutely correct makes evaluating predictions problematic. The expert annotations described by the *Drosophila* experts (Ashburner et al. 1999b) are our best available resource, but their accuracy will certainly improve as more data becomes available. At best, the data we had in hand is representative of the true situation, and our conclusions would be unchanged by using a more complete data set. At worst, there is a bias in the available data that makes our conclusions significantly misleading. We believe that the data is not unreasonably biased and that conclusions based on it are correct enough to be valuable as the basis for discussion and future development. We do not believe that the values for the various statistics introduced below are precisely what they would be using the extra information, and we emphasize that they should always be considered in the context of this particular annotated data set [for a further detailed discussion of evaluating these predictions, see Birney and Durbin (2000)].

In the next section we describe the target genomic sequence and the auxiliary data, including a critical discussion of our standard sets. Methods gives a short summary of our standard sets. Methods gives a short

Table 1. Participating Groups and Associated Annotation Categories

Program name	Program	Gene finding	Promoter recognition	EST/cDNA alignment	Protein similarity	Repeat	Gene function
Mural et al. Oakhridge, US Parra et al. Barcelona, ES	GRAIL	x	x				x
Krogh Copenhagen, DK Henikoff et al. Seattle, US Solovyev et al. Sanger, UK Gasterland et al. Rockefeller, US Benson et al. Mount Sinai, US Werner et al. Munich, GER Oliver et al. Birney, UK Reese et al. Berkeley/Santa Cruz, US	GeneID	x					x
	HMGene	x					x
	BLASTRS						x
	FGenes	x					x
	MCGIE	x	x	x	x	x	x
	TRE					x	
	CoreInspector	x					
	MCInspector	x					
	GeneWise						x
	GeneIE	x	x	x	x	x	x

10:493-501 ©2000 by Cold Spring Harbor Laboratory Press ISSN 1085-9057/00 \$5.00; www.genome.org
E-MAIL: imprese@lbl.gov; FAX (510) 486-6798.

Genome Research 48:3
www.genome.org

484 Genome Research
www.genome.org

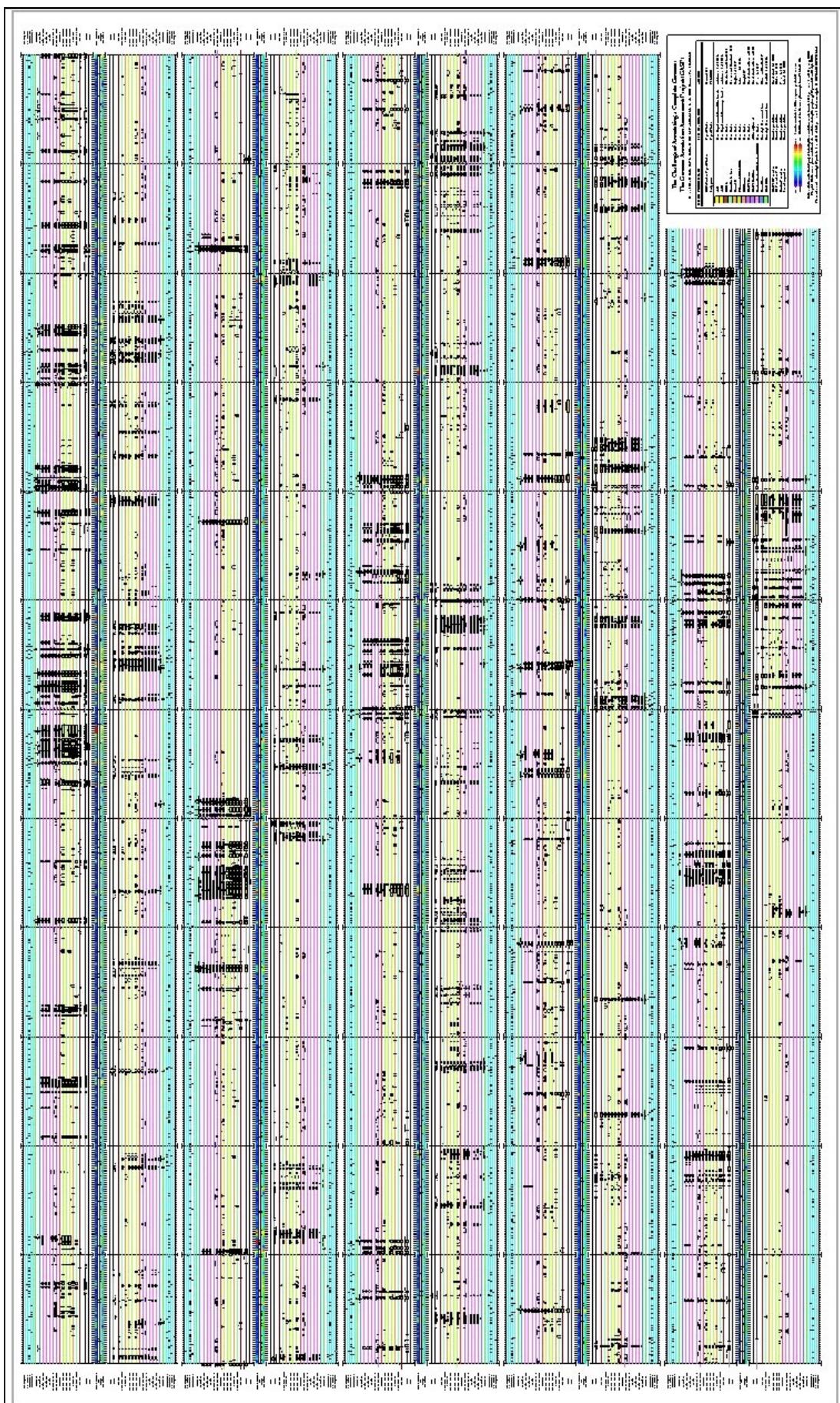


Figura 1: “*Drosophila* Genome Annotation Assessment Project”.

La mida real d'aquesta figura és de 100x60cm.

THE DROSOPHILA GENOME

REVIEW

The Genome Sequence of *Drosophila melanogaster*

Mark D. Adams,^{1*} Susan E. Antanaitis,² Robert A. Holt,¹ Cheryl A. Evans,¹ Jeannine D. Gocayne,¹ Peter G. Armanatus,¹ Steven E. Scherer,¹ Richard F. Galle,² Reed A. George,² Suzanne E. Lewis,⁴ Michael Ashburner,⁵ Scott N. Henderson,¹ Granger G. Sutton,¹ Jennifer R. Wortman,¹ Mark D. Yandell,¹ Qing Zhang,¹ Lin X. Chen,¹ Rhonda C. Brandon,¹ Yu-hui C. Rogers,¹ Robert C. Blazaj,² Mark J. Chapple,² Barret P. Beffler,² Kenneth H. Wan,² Clare M. Barton,² Gregg Hett,⁶ Catherine R. Knapp,¹ George L. Gabor Miklos,¹ Joseph F. Abril,¹ Anna Agbayani,¹ Cynthia Andrews-Pfannkoch,¹ Danita Baldwin,¹ Anand Basu,¹ James Bayendre,¹ Leyla Bayraktaroglu,¹ Ellen M. Beasley,¹ Karen Y. Besson,¹ P. V. Benos,¹ Benjamin P. Bernman,¹ Kenneth C. Burts,¹ Dana Burkhardt,¹² Michael R. Botchan,⁴ Peter Brokstein,¹ Michael C. Burstein,¹ Dana A. Busam,¹ Heather Butler,¹⁶ Edward Cardozo,¹⁷ Angela Carter,¹ Ishwar Chandra,¹ J. Michael Cherry,¹⁸ Simon Cowley,¹⁹ Carl Dahlke,¹ Lionel B. Davenport,¹ Peter Davies,¹ Beatriz de Pablo,²⁰ Arthur Deicher,¹ Michael Dzomow,²¹ Anne Dugan-Rocha,²² Shannon Dunkov,²² Patrick Dunn,¹ Kenneth J. Durbin,¹ Carlos C. Evangelista,¹ Conception Ferraz,²² Steven Ferreria,⁵ Carl Foster,¹ Andrei E. Gabrilovich,¹ Neha S. Garg,¹ William M. Gilbert,¹ Ken Glassner,¹ Anna Glodek,¹ Fangcheng Gong,¹ J. Harley Correll,¹ Zhiping Gu,¹ Ring Gaan,¹ Michael Harris,¹ Nomi L. Hirshitz,² Damon Hostin,¹ Kathryn Howland,¹ Judith R. Hernandez,³ Jarrett Houck,¹ Thomas J. Homan,¹ Timothy J. Howland,¹ Ming-hui Wei,¹ Chiyori Ibegwam,¹ Meena Jalali,¹ Francis Kalush,¹ Gary H. Karpen,¹ Zhouqi Ke,¹ James A. Kenison,²⁴ Karen A. Ketchum,¹ Bruce E. Kimmel,²⁵ Chinnappa D. Kodira,¹ Saul Kravitz,¹ David Kulp,⁶ Zhongwei Lai,¹ Paul Lasko,²⁵ Yiding Lei,¹ Alexander A. Levitsky,¹ Jaylin Li,¹ Yong Liang,¹ Xiaoying Lin,²⁶ Xianglian Liu,¹ Bettina Matthei,¹ Tina C. McIntosh,¹ Michael P. McLeod,³ Duncan McPherson,¹ Genadyi Merkulov,¹ Natalia V. Milashina,¹ Clark M. Mobarry,¹ Joe Morris,¹ Ali Mostafiz,¹ Stephen M. Mouni,²⁷ Mes Moy,¹ Brian Murphy,¹ Lee Murphy,¹ Donna M. Mousavi,¹ Michael Nelson,³ David R. Nelson,¹ Katherine Nixon,¹ Deborah R. Nusskern,¹ Joanne M. Pacile,¹ Michael Palazzo,² Gangie S. Pittman,¹ Sue Pan,¹ John Pollard,¹ Vinita Puri,¹ Martin G. Reese,¹ Knut Reinetz,¹ Karin Remington,¹ Robert D. C. Saunders,³⁰ Frederick Scheeler,¹ Hua Shen,³ Xikian Christopher Shue,¹ Inga Sidi-Klanno,¹ Michael Simpson,¹ Marian Slepnev,¹ Tom Smith,¹ Eugene Spier,¹ Allan C. Spilling,³¹ Mark Steane,¹ Renée Strang,¹ Eric Sun,¹ Robert Surkski,³² Cyndee Tector,¹ Russell Turner,¹ Elii Venter,¹ Athul H. Wang,¹ Xin Wang,¹ Zhen-Yuan Wang,¹ Trevor Woodge,¹ David A. Wassarman,³³ George M. Weinstock,¹ Jean Weissensee,¹ Sherita M. Williams,¹ Michael Williams,¹ Jane Yeh,¹ Jaydreas S. Zaveri,¹ Song Yang,² Q. Alison Yao,¹ Jane Ye,¹ Ru-Fang Yeh,¹ Ming-Zhan,¹ Guangen Zhang,¹ Qi Zhao,¹ Liansheng Zheng,¹ Xiangguo H. Zheng,¹ Fei N. Zhong,¹ Wenyan Zhang,¹ Xiaojun Zhou,¹ Shilao Gong,¹ Zhiqiang Zhu,¹ Xiaobong Zhu,¹ Hamilton O. Smith,¹ Richard A. Gibbs,³ Eugene W. Myers,¹ Gerald M. Rubin,¹ Craig Ventor¹

The fly *Drosophila melanogaster* is one of the most intensively studied organisms in biology and serves as a model system for the investigation of many developmental and cellular processes common to higher eukaryotes, including humans. We have determined the nucleotide sequence of nearly all of the ~120-megabase euchromatic portion of the *Drosophila* genome using a whole-genome shotgun sequencing strategy supported by extensive clone-based sequencing, a high-quality bacterial artificial chromosome map. Efforts are under way to close the remaining gaps in size, a third of which is centric heterochromatin. The 120 Mb of euchromatin is on two large autosomes and the X chromosome; the small fourth chromosome contains only ~1 Mb of euchromatin. The heterochromatin mainly consists of short, simple sequence elements repeated for many megabases, occasionally interrupted by inserted transposable elements, and tandem arrays of ribosomal RNA genes. It is known that there are small islands of unique sequence embedded within heterochromatin—for example, the mitogen-activated protein kinase gene *rolled* on chromosome 2, which is flanked on each side by at least 3 Mb of heterochromatin. Unlike the *C. elegans* genome, which can be completely cloned in yeast, artificial chromosomes (YACs), the simple sequence repeats are not stable over many years by the fly research community.

which has molecularly characterized genes; thus work in turn has supported the foundation for a century of genetics (⁵). Since *Drosophila* was chosen in 1906 as one of the model organisms to be studied under the auspices of the federally funded Human Genome Project, genome projects in the United States, Europe, and Canada have produced a battery of genome-wide resources (Table 1). The Berkeley and European *Drosophila* Genome Projects (BDGP and EDGP) initiated genomic sequencing (Tables 1 to 3) and finished 29 Mb of the bacterial artificial chromosome (YACs) (6) or other large-insert cloning sys-

THE DROSOPHILA GENOME

tems. This has led to a functional definition of the euchromatic genome as that portion of the genome that can be cloned stably in BACs. The euchromatic portion of the genome is the subject of both the federally funded *Drosophila* sequencing project and the work presented here. We began WGS

sequencing of *Drosophila* less than 1 year ago, with two major goals: (i) to test the strategy on a large and complex eukaryotic genome as a prelude to sequencing the human genome, and (ii) to provide a complete, high-quality genomic sequence to the *Drosophila* research community so as to advance research in this important model organism.

WGS sequencing is an effective and efficient way to sequence the genomes of organisms, which are generally between 0.5 and 6 Mb in size (7). In this strategy, all the DNA of an organism is sheared into segments a few thousand base pairs (bp) in length and cloned directly into a plasmid vector suitable for DNA sequencing. Sufficient DNA is covered numerous times, in fragments of ~500 bp. After sequencing, the fragments are assembled in overlapping segments to reconstruct the complete genome sequence.

In addition to their much larger size, eukaryotic genomes often contain substantial amounts of repetitive sequence that can be potential to interfere with correct sequence assembly. Weber and Myers (8) presented a theoretical analysis of WGS sequencing, in which they examined the impact of repetitive sequences, discussed experimental strategies to mitigate their effect on sequence assembly, and suggested that the WGS method could be applied effectively to large eukaryotic genomes. A key component of the strategy is obtaining sequence data from each end of the cloned ends—sequences (mate pairs) is a critical element in producing a correct assembly.

Genomic Structure

WGS libraries were prepared with three different insert sizes. The 14-kb clones are large enough to span the most common repetitive sequence elements in *Drosophila*, the retrotransposons. End-sequence data from the BACs provided long-range linking information that was used to confirm the overall structure of the assembly (9). More than 3 million sequence reads were ob-

tained from whole-genome libraries (Fig. 2) and contained heterochromatic simple sequence repeats, indicating that the heterochromatic DNA is not stably cloned in the small-insert vectors used for the WGS libraries. A BAC-based physical map spanning >95% of the euchromatic portion of the genome was constructed by screening a BAC library with sequence-tagged site (STS) markers (¹⁰). More than 29 Mb of high-quality finished sequence have been completed from BAC, P1, and cosmid clones, and draft sequence data (~1.5× average coverage) were obtained from an additional 82 BAC and P1 clones spanning in total >90% of the genome (Table 3). The clone-based draft sequence served two purposes. It improved the likelihood of accurate assembly, and it allowed the identification of templates and primers for filling gaps that remain after assembly. An initial assembly was performed using the WGS data and BAC end-sequence (WGS-only) assembly (4); subsequent assemblies used the clone-based draft sequence data (Joint assembly; Fig. 3 and Table 3) illustrate the status of the euchromatic sequence resulting from each of these assemblies and the current status following the directed gap closure completed to date. The sequence assembly is described in detail in an accompanying paper (11).

Assembly resulted in a set of "scaffolds."

Each scaffold is a set of contiguous sequences (contigs) ordered and oriented with respect to one another by matepairs such that the gaps between adjacent contigs are of known size and are spanned by clones with end-sequences flanking the gap. Gaps within scaffolds are called sequence gaps; gaps between scaffolds are called "physical gaps" because there are no clones identified spanning the gap. Two methods were used to map the scaffolds to chromosomes: (i) cross-referencing between STS numbers assigned in the assembled sequence and the BAC-based STS content map, and (ii) cross-referencing between assembled sequence data and shotgun sequence data obtained from individual long-patch clones selected from the BAC physical map. The mapped scaffolds from the joint assembly, totaling 116.2 Mb after initial

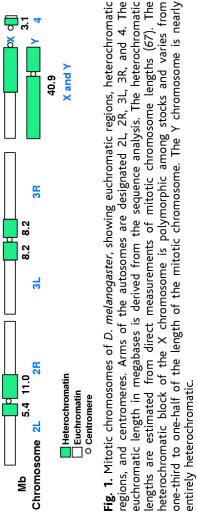


FIG. 1. Minichromosomes of *D. melanogaster*, showing euchromatic regions, heterochromatic regions, and centromeres. Arms of the autosomes are designated 2L, 2R, 3L, and 4. The euchromatic length in megabases is derived from the sequence analysis. The heterochromatic lengths are estimated from direct measurements of mitotic chromosome lengths (67). The heterochromatic block of the X chromosome is polymorphic among stocks and values from one-third to one-half of the length of the mitotic chromosome. The Y chromosome is nearly entirely heterochromatic.

*To whom correspondence should be addressed.

24 MARCH 2000 VOL 287 SCIENCE www.sciencemag.org

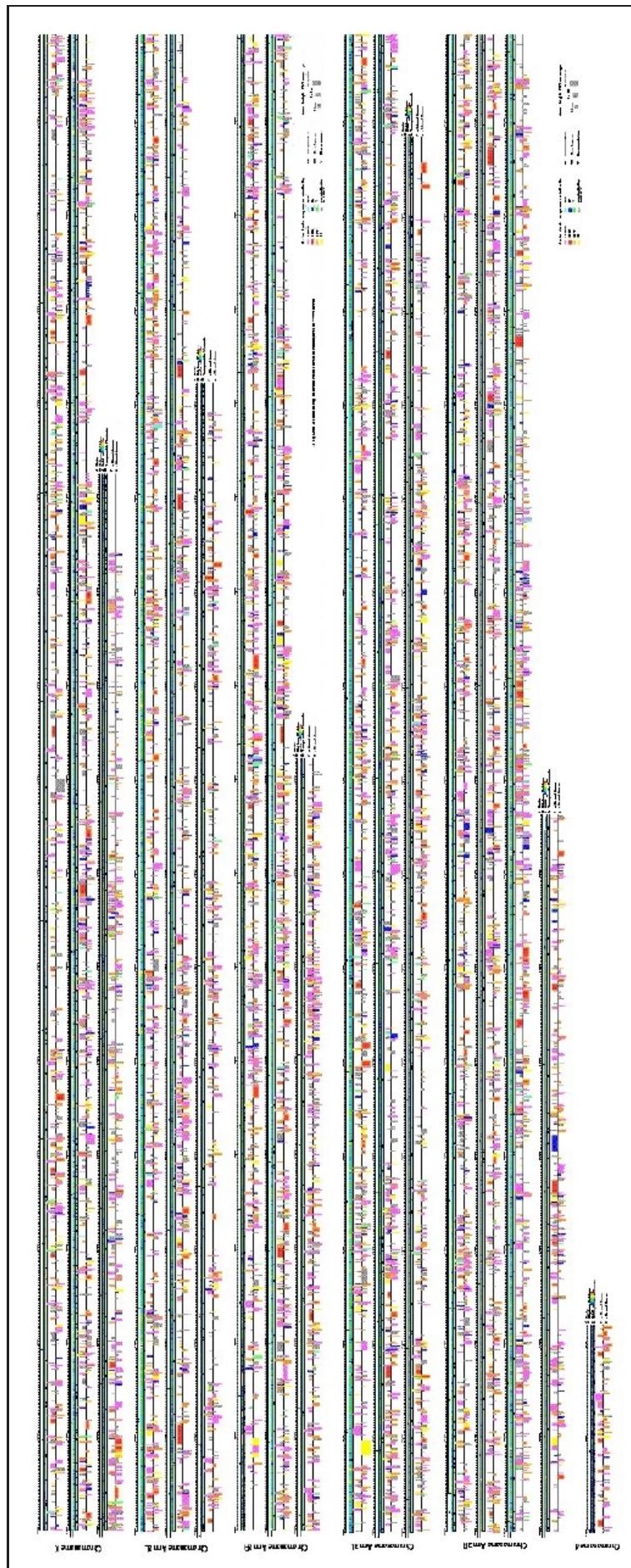


Figura 2: “Coding Content of the *Drosophila Genome*”.

La mida real d'aquesta figura és de dues pàgines de 108x28cm.

THE HUMAN GENOME

The Sequence of the Human Genome

- J. Craig Venter,^{1,*} Mark D. Adams,¹ Elgene W. Myers,¹ Peter W. Li,¹ Richard J. Mural,¹ Granger G. Sutton,¹ Hamilton O. Smith,¹ Mark Fandell,¹ Cheryl A. Evans,¹ Robert A. Holt,¹ Jeannine D. Gocayne,¹ Peter Amanatidis,¹ Richard M. Balow,¹ Daniel H. Huson,¹ Jennifer Russo-Wortman,¹ Qing Zhang,¹ Chinnappan K., Xiangqun H. Zheng,¹ Lin Chen,¹ Marian Supaski,¹ Gangadarhan Subramanian,¹ Paul D. Thomas,¹ Jinghui Zhang,¹ George L. Gibson Miklos,² Catherine Nelson³ Samuel Broder,¹ Andrew G. Clark,⁴ Joe Nadeau,⁵ Victor A. McSwick,⁶ Norton Zinder,⁷ Arnold J. Levine,⁷ Richard A. Roberts,⁸ Mel Simon,⁹ Carolyn Stayman,¹⁰ Michael Hunkapiller,¹ Randall Bolanos,¹ Arthur Dancher,¹ Ian Dew,¹ Daniel Fastuol,¹ Michael Flanigan,¹ Lilliana Flores,¹ Aaron Halpern,¹ Sridhar Hannenhalli,¹ Saul Kravitz,¹ Samuel Levy,¹ Clark Moberly,¹ Kent Reiner,¹ Karin Ramington,¹ Jane Abu-Threideh,¹ Ellen Brasley,¹ Kendra Biddick,¹ Vivien Bonazzi,¹ Rhonda Brandon,¹ Michelle Cargil,¹ Ishwar Chandranoulsivaram,¹ Rosanne Charlab,¹ Kabir Chatravedi,¹ Zuoming Deng,¹ Valentina Di Francesco,¹ Patrick Dunn,¹ Karen Eilbeck,¹ George Evangelista,¹ Andrei E. Gabrilian,¹ Weinui Gan,¹ Fangcheng Gong,¹ Zhiping Gu,¹ Ping Guan,¹ Thomas J. Heiman,¹ Maureen E. Higgins,¹ Rui-Pu Ji,¹ Zhaoxi Ke,¹ Karen A. Ketcham,¹ Zhongwu Lai,¹ Yiding Lei,¹ Zhengyu Li,¹ Yong Liang,¹ Xiaoying Lin,¹ Fu Lu,¹ Helen M. Moore,¹ Ashwinikumar K. Naik,¹ Gemma V. Merkulov,¹ Natalia Mishina,¹ Beena Neelam,¹ Deborah Nuesslein,¹ Douglas B. Rusch,¹ Steven Sanberg,¹² Vaibhav A. Narayan,¹ Shaoqiang C. Zhu,¹ Shaying Zhao,¹ Jianhua Yan,¹ Alison Yao,¹ Xin Wang,¹ Jian Wang,¹ Wei Shio,¹ Bixiong Shue,¹ Jingtang Sun,¹ Zhen Yuan Wang,¹ Alfonso Yee,¹ Ming Zhang,¹ Weiqing Zhang,¹ Hongyu Zhang,¹ Qi Zhao,¹ Liangsheng Zheng,¹ Fei Zhou,¹ Wenyan Zhong,¹ Shaoping C. Zhu,¹ Shaying Zhao,¹ Dennis Gilbert,¹ Suzanne Baumhueter,¹ Gene Spier,¹ Christine Carter,¹ Anibal Craychik,¹ Trevor Woodage,¹ Feroze Ali,¹ Hujin An,¹ Adenore Aye,¹ Danita Baldwin,¹ Holly Baden,¹ Mary Barnstead,¹ Ian Barrow,¹ Karen Beeson,¹ Dara Busam,¹ Amy Carver,¹ Angela Center,¹ Ming Lai Cheng,¹ Liz Curry,¹ Steve Danaher,¹ Lionel Daverden,¹ Raymond Desilets,¹ Susanne Dietz,¹ Kristina Dooson,¹ Lisa Douc,¹ Steven Ferrera,¹ Neha Garg,¹ Andres Gluecksmann,¹ Britt Hart,¹ Jason Haynes,¹ Charles Haynes,¹ Jeffrey Heiner,¹ Suzanne Hladun,¹ Damon Hostin,¹ Jarrett Houck,¹ Timothy Howland,¹ Chinvery Isbergwan,¹ Jeffery Johnson,¹ Francis Kalush,¹ Leahy Kline,¹ Shashi Kaduru,¹ Amy Love,¹ Felicia Mann,¹ David May,¹ Steven McCawley,¹ Tina McIntosh,¹ Ivy McMillen,¹ Mee Moy,¹ Linda Moy,¹ Brian Murphy,¹ Keith Nelson,¹ Cynthia Pfannkoch,¹ Eric Pratts,¹ Vinita Pur,¹ Hina Qureshi,¹ Matthew Reardon,¹ Robert Rodriguez,¹ Yu-Hui Rogers,¹ Deanna Rombola,¹ Bob Rutledge,¹ Richard Scott,¹ Cynthia Stitter,¹ Michelle Smalikow,¹ Erin Stewart,¹ René Strong,¹ Ellen Suh,¹ Reginald Thomas,¹ James Tint,¹ Sukyee Tse,¹ Claire Vech,¹ Gary Wang,¹ Jeremy Wetter,¹ Sherita Williams,¹ Monica Williams,¹ Sandra Windsor,¹ Emily Winn-Den,¹ Kerilen Wolfe,¹ Jayshree Zaveri,¹ Karena Zaveri,¹ Joseph F. Abri,^{1,*} Roderic Guig,¹⁴ Michael J. Campbell,¹ Kimmie L. Campbell,¹ Brian Karlik,¹ Anish Kaparaju,¹ Huayui Mi,¹ Betty Lazareva,¹ Thomas Hatton,¹ Apurva Narechania,¹ Karen Diemer,¹ Anushya Muruganujan,¹ Nan Guo,¹ Shinji Kashi,¹ Ross Lipert,¹ Russell Schwartz,¹ Brian Walenz,¹ Marcelo Caminha,¹ John Carnes-Stine,¹ Parris Caull,¹ James Bayendale,¹ Louis Bick,¹ Carl Dahke,¹ Anne Deslatte Mays,¹ Maria Dambroski,¹ Michael Donnelly,¹ Stephen Esparham,¹ Carl Fosler,¹ Harold Gire,¹ Kenneth Glasser,¹ Anna Glodek,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 Jordam,¹ John Kasha,¹ Leonid Kogan,¹ Cheryl Kraft,¹ Alexander Levitsky,¹ Mark Lewis,¹ John Lopez,¹ Daniel Ma,¹ William Maioros,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 Mihayuan Wen,¹ David Wu,¹ Mitchell Wu,¹ Mitchel Wu,¹ Ashton Xia,¹ Ali Zandieh,¹ Xiaohong Zhu¹

THE HUMAN GENOME

A 2.91-billion base pair (bp) consensus sequence of the euchromatic portion of the human genome was generated by the whole-genome shotgun sequencing method. The 1.46-billion bp DNA sequence was generated over 9 months from 22,271,853 high-quality sequence reads (5.11-fold coverage of the genome) from both ends of plasmid clones made from the DNA of five individuals. Two assembly strategies—a whole-genome assembly and a regional chromosome assembly—were used, each combining sequence data from Celera and the publicly funded genome effort. The public data were shredded into 550-bp sequence fragments to create a 29-fold coverage of those genome regions that had been sequenced, without including biases inherent in the cloning and assembly procedure used by the publicly funded group. This brought the effective coverage in the assemblies to 5.11 times the number and size of gaps in the final assembly over what would be obtained with 5.11-fold coverage. The two assembly strategies yielded very similar results that largely agree with independent mapping data. The assemblies effectively cover the euchromatic portion of the human chromosomes. More than 90% of the genome is in scaffold assemblies of 100,000 bp or more, and 25% of the genome is in scaffolds of 10 million bp or larger. Analysis of the genome sequence revealed 2,588 protein-encoding transcripts for which there was strong corroborating evidence and an additional ~12,000 computationally derived genes with mouse orthologs or other weak supporting evidence. Although gene-density clusters are obvious, almost half the genome is dispersed in low G+C sequences separated by large tracts of apparently noncoding sequence. Only 1% of the genome is spanned by exons, whereas 24% is in introns, with 75% of the genome being intergenic DNA. Duplications of segmental blocks, ranging in size up to chromosomal lengths, are abundant throughout the genome and reveal a complex evolutionary history. Comparative genomic analysis indicates vertebrate exon density and distribution with neuronal function, with tissue-specific developmental regulation, and with the hemostatic and immune systems. DNA sequence comparisons between the consensus sequence and publicly funded genome data provided locations of 2.1 million single-nucleotide polymorphisms (SNPs). A random pair of human haploid genomes differed at a rate of 1 bp per 1250 on average, but there was marked heterogeneity in the level of polymorphism across the genome. Less than 1% of all SNPs resulted in variation in genome sequence, but the task of determining which SNPs have functional consequences remains an open challenge.

Decoding of the DNA that constitutes the human genome has been widely anticipated for the contribution it will make toward understanding human evolution, the causation of disease, and the interplay between the environment and heredity in defining the human condition. A project with the goal of determining the complete nucleotide sequence of the human genome was first undertaken in 1985 (¹). In subsequent years, the idea met with mixed reactions in the scientific community. However, in 1990, the Human Genome Project (HGP) was officially initiated in the United States under the direction of the National Institutes of Health. In 1992, the Wellcome Trust, the Wellcome Trust Sanger Institute, and the Max Planck Institute for Molecular Genetics announced their intent to build a unique genome-sequencing facility to determine the sequence of the human genome over a 3-year period. Here we report the penultimate milestone along the path toward that goal, a nearly complete sequence of the euchromatic portion of the human genome. The sequencing of the genome sequence. In 1998 we announced our intention to build a unique genome-sequencing facility to determine the sequence of the human genome over a 3-year period. Here we report the penultimate milestone along the path toward that goal, a nearly complete sequence of the euchromatic portion of the human genome. The sequencing of the genome was performed by a whole-genome shotgun method with subsequent assembly of the sequenced segments. The modern history of DNA sequencing began in 1977, when Sanger reported his method for determining the order of nucleotides of DNA using chain-terminating nucleotide analogs (³). In the same year, the first human gene was isolated and sequenced (⁴). In 1986, Hood and co-workers (⁵) described an improvement in the Sanger sequencing method that included attaching fluorescent dyes to the nucleotides which permitted them to be sequentially read by a computer. The first automated DNA sequencer, developed by Applied Biosystems in California in 1987, was shown to be successful when the sequences of two genes were obtained with this new technology (⁶). From early sequencing of human genomic regions (⁷), it became clear that cDNA sequences (which are reverse-transcribed from RNA) would be essential to annotate and validate gene predictions in the human genome. These studies were the basis in part for the development of the expressed sequence tag (EST) method of gene identification (⁸), which is a random selection, very high throughput sequencing approach to characterize cDNA libraries. The EST method led to the rapid discovery and mapping of human genes (⁹). The increasing numbers of human EST sequences necessitated the development of new computer algorithms to analyze large amounts of sequence data and, in 1993, at the Institute for Genomic Research (TIGR), an algorithm was developed that permitted assembly and analysis of hundreds of thousands of ESTs. This algorithm permitted characterization and annotation of human genes on the basis of 30,000 EST assemblies (¹⁰). The complete 49-kbp bacteriophage lambda genome sequence was determined by a shotgun restriction digest method in 1982 (¹¹). When considering methods for sequencing the simiploidy virus genome in 1991 (¹²), a whole-genome shotgun sequencing method was discussed and subsequently rejected owing to the lack of appropriate software tools for genome assembly. However, in 1994, when a microbial genome-sequencing project was undertaken at TIGR, a whole-genome shotgun sequencing approach was considered possible with the IGR EST assembly algorithm (¹³). In 1995, the 1.8-Mbp *Haemophilus influenzae* genome was completed by a whole-genome shotgun sequencing method (¹⁴). The experience with several subsequent genome-sequencing efforts established the broad applicability of this approach (¹⁴,¹⁵). A key feature of the sequencing approach used for these megabase-size and larger genomes was the use of paired-end sequences (also called mate pairs), derived from subclone libraries with distinct insert sizes and cloning characteristics. Paired-end sequences are sequences 500 to 600 bp in length from both ends of double-stranded DNA clones of prescribed lengths. The success of end sequencing from long segments (18 to 20 kbp) of DNA cloned into bacteriophage lambda in assembly of the microbial genomes led to the suggestion (¹⁶) of an approach to simultaneously determine the order of nucleotides of an genome (¹⁶).

1305

www.science.org SCIENCE VOL 291 16 FEBRUARY 2001

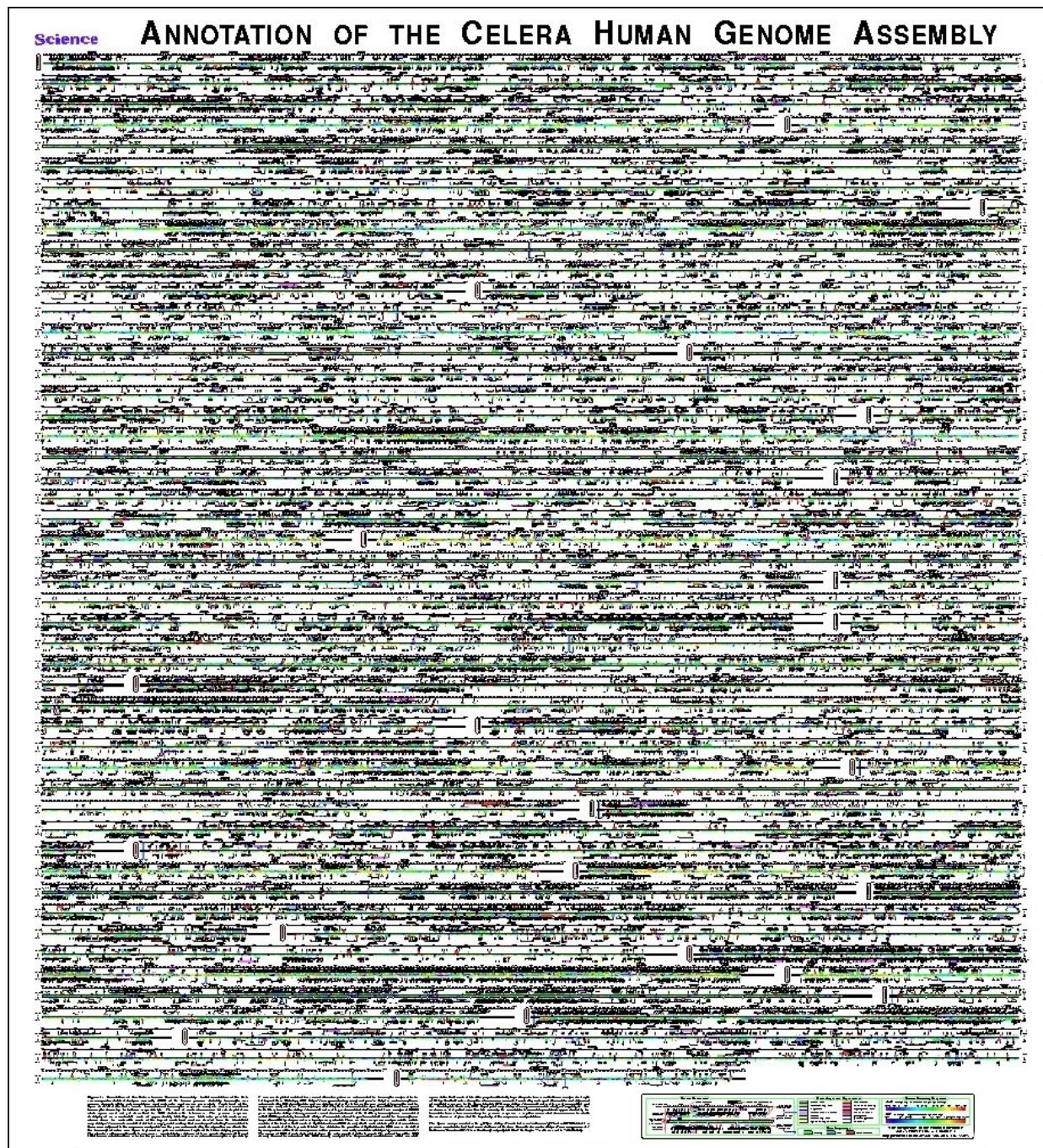


Figura 3: “Annotation of the Celera Human Genome Assembly”.

La mida real d'aquesta figura és de 106x150cm.

Methods

An Assessment of Gene Prediction Accuracy in Large DNA Sequences

Roderic Guigó,^{1,3} Pankaj Agarwal,² Josep F. Abril,¹ Moisés Burset,¹ and James W. Fickett.²

¹Grup de Recerca en Informàtica Mèdica, Institut Municipal d'Investigació Mèdica, Universitat Pompeu Fabra, E-08003 Barcelona, Spain; ²Department of Bioinformatics, SmithKline Beecham Pharmaceuticals Research and Development, King of Prussia, Pennsylvania 19406, USA

One of the first useful products from the human genome will be a set of predicted genes. Besides its intrinsic scientific interest, the accuracy and completeness of this data set is of considerable importance for human health and medicine. Though progress has been made on computational gene identification in terms of both methods and accuracy evaluation measures, most of the sequence sets in which the programs are tested are short genomic sequences, and there is concern that these accuracy measures may not extrapolate well to larger, more challenging data sets. Given the absence of experimentally verified large genomic data sets, we constructed a semiartificial test set comprising a number of short single-gene genomic sequences with randomly generated intervening regions. This test set, which should still present an easier problem than real human genomic sequence, mimics the ~200kb long BACs being sequenced. In our experiments with these longer genomic sequences, the accuracy of GENSCAN, one of the most accurate ab initio gene prediction programs, dropped significantly, although its sensitivity remained high. Conversely, the accuracy of similarity-based programs, such as GENEFEST, PROCUSTES, and BLASTX, was not affected significantly by the presence of random intervening sequence, but depended on the strength of the similarity to the protein homolog. As expected, the accuracy dropped if the models were built using more distant homologs, and we were able to quantitatively estimate this decline. However, the specificities of these techniques are still rather good even when the similarity is weak, which is a desirable characteristic for driving expensive follow-up experiments. Our experiments suggest that though gene prediction will improve with every new protein that is discovered and through improvements in the current set of tools, we still have a long way to go before we can decipher the precise exonic structure of every gene in the human genome using purely computational methodology.

The nucleotide genomic sequence is the primary product of the Human Genome Project, but a major short- and mid-term interest will be the amino acid sequences of the proteins encoded in the genome. Thus, methods that reliably predict the genes encoded in genomic sequence are essential, and computational gene identification continues to be an active field of research (for reviews, see Fickett 1996; Claverie 1997; Guigó 1997a; Burge and Karlin 1998; Haussler 1998). A new generation of gene prediction programs based on Hidden Markov Models (Burge and Karlin 1997) have shown significantly greater accuracy than previous programs based on other methodologies (Burset and Guigó 1996). Conversely, as the databases of known coding sequences increase in size, gene prediction methods based on sequence similarity to coding sequences, mainly proteins and ESTs, are becoming increasingly useful and are routinely used to identify putative genes in genomic sequences (The *C. elegans* Sequencing Consortium 1998). We have recently published an evaluation of gene prediction programs based on HMMs in genomic sequences (The *C. elegans* Sequencing Consortium 1998). We have recently published an evaluation of gene prediction programs based on HMMs in genomic sequences (The *C. elegans* Sequencing Consortium 1998). We have recently published an evaluation of gene prediction programs based on HMMs in genomic sequences (The *C. elegans* Sequencing Consortium 1998).

*Corresponding author.

E-mail: rguigo@imim.es; FAX 34931-221-3237.

Article and publication are at www.genome.org/cgi/do/10.1093/gg/122800.

Guigó et al.

EST sequences, such as EST_GENOME (Mott 1997), could also be included in this category. These programs promise highly accurate predictions, but at the cost of greater computational time. However, this increase in accuracy has not been well-quantified on challenging data sets. The effects of the degree of similarity between the candidate homolog and the genomic sequence also deserve careful evaluation.

We believe a more realistic evaluation of the currently available gene prediction tools on challenging data sets would be useful. Ideally, one would like to benchmark the computational gene identification programs in real genomic sequences. The main problem is that most real sequences the structure of the genes has not been verified exhaustively by experimental means, and thus it is impossible to calibrate the accuracy of the predictions. Only recently, extensively annotated large genomic sequences from higher eukaryotic organisms have become available from the human genome (<http://www.igap.mncacub/GeneSet>) and from the fly genome (<http://www.fruitfly.org/GASPI>). In spite of the experimental analysis, the possibility of understanding genes in the sequence cannot be easily ruled out, which makes accuracy difficult to measure. Here, we attempt to overcome the lack of well-annotated large genomic sequences by constructing semiartificial ones. In these semiartificial sequences, known genomic sequences have been embedded in simulated integenic DNA, and therefore, the location of all coding exons is known. Although the approach may seem unrealistic, we believe that the results obtained are instructive with regard to the accuracy of currently available gene identification tools.

We evaluate the accuracy of representatives of a wide variety of computational gene identification approaches: GENSCAN (Burge and Karlin 1997), an ab initio geneidifier; BLASTX (Altschul et al. 1990; Gish and States 1993), a genefinding-oriented similarity search program; and PROCUSTES (Gelfand et al. 1996) and GENEFEST (Birney and Durbin 1997), genefinders based on aligning genomic DNA sequence fragment to a homologous protein sequence. We evaluate these programs on two benchmark data sets: A set of well-

annotated single-gene DNA sequences, and a set of semiartificial genomic (SAG) sequences created by embedding the single-gene sequences from the first data set in simulated intergenic DNA.

RESULTS

We investigated the accuracy of the gene prediction tools (GENSCAN, PROCUSTES, GENEFEST, SAG, BLASTX) described in Methods on two benchmark sets. In all cases, sequences were masked previously for repeated regions using EXPANDMASKER (A. Shatt and P. Green, unpubl.). The gene predictions obtained using the different tools were compared with the actual gene annotations using the accuracy measures described in Methods.

Accuracy in Single Gene Sequences

Table 1 shows the accuracy of the different gene prediction tools on h178, the set of single gene sequences. GENSCAN's accuracy is comparable to that reported earlier (Burge and Karlin 1997). On average, 90% of the coding nucleotides and 70% of the exons are predicted correctly by GENSCAN. Only 7% of the actual exons are missed completely, and only 9% of the predicted exons are wrong. We believe this is close to the maximum accuracy that can be achieved using currently available ab initio gene prediction programs.

The quality of the gene models inferred from BLASTX depends on the strategy used. Default usage of BLASTX produced poorer predictions than more sophisticated strategies. Results for BLASTX default correspond to those published in Guigó et al. (2000). Discrepancies between numbers in Table 1 and those reported in Guigó et al. (2000) are due to the differences in the way the accuracy measures are summarized. In Guigó et al. (2000), we computed the accuracy measures on each test sequence and averaged all of them. Here, we compute the accuracy measures globally from the total number of prediction successes (at the base or exon level) on all sequences. The default BLASTX strategy produces reasonably high sensitivity (0.91) by projecting all HSPs over a given threshold along the query DNA sequence, but the sensitivity rises to an amazing 0.97, if the topcombn feature

Table 1. Accuracy of Gene Prediction Tools in the Set of Single Gene Sequences (h178)

Program	Nucleotide No.	Exon								
		Sn	Sp	CC	Sn	Sp	2	Sn + Sp	ME	WE
GenScan	177	0.93	0.90	0.78	0.75	0.76	0.08	0.10		
Blastx default	175	0.91	0.79	0.82	0.04	0.04	0.04	0.12	0.05	
Blastx topcombn	174	0.97	0.86	0.04	0.04	0.04	0.04	0.08	0.02	
Blastx 2 stages	175	0.90	0.92	0.90	0.10	0.12	0.11	0.19	0.02	
Genefest	177	0.98	0.98	0.97	0.88	0.91	0.89	0.96		
Procuress	177	0.93	0.95	0.93	0.76	0.82	0.79	0.11	0.04	