

Genome Annotation Assessment *Drosophila melanogaster*

Drosophila

Martin Reese*, George Hartzell*, Norman Harris*, Uwe Ohler*,
Lewis*

† and Suzanne E.

*Berkeley Drosophila Genome Project, Department of Molecular and Cell
Biology,
University of California, Berkeley, California 94720-3200

Biology,

†Chair of Pattern Recognition, University of Erlangen-Nuremberg,
Erlangen, Germany.

rtenss@-91058

Send to:

Martin Reese
Lawrence Berkeley National Laboratory
Cyclotron Road MS-121
Berkeley, CA 94720
Phone (510) 86-4800
Fax (510) 86-6798
E-mail: mgreese@lbl.gov

Abstract

Computational methods for automated genome annotation are critical for community ability to make full use of genomic sequence being generated and released. To explore the accuracy of automated feature prediction in the genome of higher organisms we evaluated the performance of a large, well-characterized query set from the *Drosophila melanogaster* Genome Annotation Assessment Project (GASP), launched in May 1999, by applying state-of-the-art tools contributed by a number of sites and repeat element evaluators. The predictions in this set, based on a previously released high quality full-length DNA sequence assembly, are of a high depth of study of regions of *Drosophila* experts (Ashburner *et al.* 1999b). While these standards are only approximately known distributions of features in this region, we believe that a more detailed evaluation of these features is a meaningful contribution to the field. The results were presented at the International Systemic Molecular Biology (ISMB-99) August 1999 (Reese *et al.* 1999) and 95% of the coding nucleotide in the region were correctly identified by major methods. In addition, the correct introns/exon structures were predicted for 74% of genes. Homology-based annotation techniques recognized and associated functions with most of the genes in the region that remained only identified by the *ab initio* techniques. This experiment presents the first assessment of promoter prediction techniques. A significant number of large contiguously discovered promoter predictors with high positive rates and their prediction difficulties in integrating information from DNA/EST alignments with promoter prediction. The number of false positive classifications was discovered to be less than a third of the promoter regions. We believe that establishing standards for evaluating genomic annotations is an essential part of the performance of existing automated genome annotation tools.

experiment establish baseline which contribute the value of
project and guide further research in bioinformatics.

ongoing large-scale annotation

Introduction to the Genom Annotation Assessment

Project (GASP)

Genome annotation is rapidly evolving in genomic capabilities. The generation of genomic sequences and the predominant use of computational annotation processes is significant information possible. The genome is emphasized in the location and structure of genes. This *ab initio* gene finding by identifying homologous genes on the organism alignment of full-length partial RNA sequences to the genomic of the host. Related techniques can be used to identify the of regulatory elements or repetitive sequence elements. That the functional classification of identified genes current genes with known functions.

large-scale of the complete sequence of the genome is accomplished by the DNA through combinations of features such as the image of genome annotation, depending on discovering homologies

With the objective of assessment of the accuracy of genome annotation. The GASP project is organized to encourage the development of new models for the improvement of existing approaches through the use of prediction made by current state-of-the-art programs.

that automated techniques to formulate guidelines for the development of new assessment and comparison

The GASP experiment finds and, similarly, Assessment techniques for protein structure prediction contests (Dunbrack *et al.* 1997; Levitt 1997; Moulton *et al.* 1997; Moulton *et al.* 1997; Sippl *et al.* 1999; Zemla *et al.* 1999) described at <http://predictioncenter.llnl.gov>. However, unlike the context, GASP promotes collaboration to evaluate various annotation.

ways. The CASP Critical protein structure prediction CASP techniques for genome

The ASB experiment consisted of the following stages:

- Training of the *Ad* region including megabases of *Drosophila melanogaster* genomic sequence was collected by organizers and provided to participants.
- Standards developed to evaluate submissions while participating groups produced and submitted an annotation for the region.
- The participating group predictions were compared to the standards, team independent assessors evaluated the results, and the results were presented at ISMB-99.

Participants were given finished sequences for *Ad* regions and related training that they have access to full-length DNA sequences that sequence data by Ashburner *et al.* (1999) describe the *Ad* gene. The experiment was widely announced to any participants who submitted low-cost and encouraged disclosure of their methods. Since we fortunately attract a large group of participants who provide a variety of annotations, we believe that the non-manual annotation.

Twelve groups participated in ASB, submitting annotations in 11 categories: *ab* *init* *gen* *end* *promoter* *recognition* *EST/cDNA* *alignment* *protein* *similarity* *repetitive* *sequence* *identification* *gene* *function*. Tables of participating groups and their annotations are available in the ASB submission. Additional papers have been written by participants themselves and describe their methods and results in detail.

should not be taken as a standard, but as a guide to what is possible. It is not a goal to have a complete set of annotations, but to have a set that is useful for the community. The current annotations are a good starting point, but they are not perfect. There are many gaps, and it is important to identify these gaps and work to fill them. The community should be encouraged to contribute to the improvement of the annotations, and to share their own data and methods. This will help to ensure that the annotations are accurate, complete, and up-to-date. The community should also be encouraged to use the annotations in a way that is consistent with the standards, and to report any errors or omissions. This will help to ensure that the annotations are reliable and useful for the community.

In this section, we describe the current state of the annotations, and discuss the challenges that we face. We also describe the methods that we used to generate the annotations, and the results of our analysis. We conclude by discussing the future of the annotations, and the role of the community in their development.

Data Benchmark *Drosophila melanogaster*

The selection of a target region for annotation is a difficult task, and one that is often done in an ad hoc manner. The current annotations were generated using a set of criteria that were designed to ensure that the annotations were accurate, complete, and up-to-date. The community should be encouraged to contribute to the improvement of the annotations, and to share their own data and methods. This will help to ensure that the annotations are reliable and useful for the community.

good experimentally verified correct answers should be anonymous and blind. It would be possible. The *Adh* gene of the *Drosophila melanogaster* genome meets the criteria. *Drosophila melanogaster* is the most important model organism, and though the *Adh* region has been extensively studied, the gene annotation and cDNAs for the region are not published until the conclusion of the ASP experiment. The 9 base *Adh* contigs are large enough to challenge containing genes with a variety of structures and included regions of high density. They are completely blind, and to we even in several cDNA and genomic sequence known in the region are available. The experiment.

2. Genomic DNA Sequence

The contiguous genomic sequence of the *Adh* region in the *Drosophila melanogaster* genome spans nearly 1000 bases and the sequence is a series of overlapping fragments. Clones from the Berkeley *Drosophila* Genom Project (Rubin, 1999) and the European *Drosophila* Genom Project (Ashburner, 1999). This sequence is believed to be of high quality with an estimated error rate of 10,000 bases per PHRAP quality score. A detailed analysis of this region is accessible through the DGDG website (<http://www.fruitfly.org/publications/Adh.html>) by Ashburner *et al.* (1999b).

2. Curating the Sequences

We provided several *Drosophila melanogaster* specific ASP participants. This enabled participants to help in the *Drosophila* and facilitate comparison of various approaches that are biased by organism specific files. Following rates of sequences, extracted from Flybase and FMBI provided by the European *Drosophila* Genom Project.

Cambridge and the Berkeley Drosophila Genome Project were made available in

found at <http://www.fruitfly.org/GASP/data/data.html>:

- 36 complete coding sequences (start codon), excluding transposable elements, pseudogenes, non-coding RNAs, mitochondrial DNA sequences (3,123 entries);
- Non-redundant repetitive sequences including transposable elements (96 entries);
- Transposon sequences containing only the longest sequence of each transposon family, excluding defective transposable elements (44 entries);
- Genomic DNA from 275 multiple single-exon non-redundant genes, other than those from the multiple transposon families taken from GenBank version 109;
- 366 related promoter regions taken from EPD (Cavérier *et al.* 1990, Cavérier *et al.* 2000) and collection made by Arkhipova (1995);
- An uncurated set of DNA ATSS sequences from a progress of the Berkeley Drosophila Genome Project.

Five of the participating groups reported making of these datasets.

2.3. Resources for assessing prediction of the correct answer

Comparative studies have been used to evaluate the most important factor in determining the usefulness of the study results. For meaningful standards to be appropriate to the study's goals, the study's goals should be to evaluate tools that predict gene structure in complex eukaryotic organisms. We use a standard

complex eukaryotic genomes, choosing to use the megabase sequence from the *Adelphi* of *Drosophila melanogaster*. Comparing predicted annotations in regions only on sequential standards, we believe that the correct ones are established by techniques that are independent of the approaches being studied. The prediction of prior knowledge of standards, which would contain the correct structural genes in regions without extraneous annotations, is unfortunately not possible to obtain the underlying biology completely understood. What is a two-part approximation of the perfect data, taking advantage of the RGD DNA sequencing project (<http://www.fruitfly.org/EST>), the Drosophila community forum, and annotated regions (Ashburner *et al.* 1999). One component of the *std* data sets is high quality sequence of full-length DNA, one of the *Adelphi* providing a standard with annotations that are likely to be uncertainly not exhaustive. The second component of the *std* data set is the annotations being developed for Ashburner *et al.* (1999) by standard with complete coverage of regions, although this confidence in the accuracy and independence of the annotations is believed that this two-part approximation allows a useful conclusion about the ability to accurately predict gene structure in complex eukaryotic organisms, even though the absolutely perfect data do not exist.

Eukaryotic transcription annotations have complex structures as the composition of fundamental features such as the 5' cap, the transcription start site, the transcription start site (TSS), start of primary and principal boundaries, and the poly-adenylation signal, start positions, and coding start sites. One prediction evaluation focused on annotations has specific coding regions, the start codon through various intron-exon boundaries, the coding region, and the annotations. While the features are biologically interesting, we can make a reliable

method for evaluating their predictions. Whenever possible, we
evidence for our evaluations, where not available, combine
curated domain experts.

an unambiguous biological
several types of evidence

Of our first standard, called *stdv*, to build a set of annotations we believed
were likely to be correct in the future (e.g.,
we were able to include divergent regions based
full-length DNA sequences from this region with high quality
The DNA sequences are the product of a DNA sequencing project
Drosophila Genome Project and have been submitted to GenBank
Working for the DNA libraries, the longest for each gene
sequencing high quality. Starting with these DNA sequences
the genomic sequences in *mm4* (Florea *et al.* 1998) and filter them several times
eighty candidate DNA sequences, three paralogous genes
appeared as a non-coding artifact (unspliced RNA) in multiple inserts
with alignment for fifty-eight DNA clones. These alignments
splicing. We require that all proposed splice
coordinates for splice sites respectively and that they correspond to high quality
threshold high 98% positive rate for splice sites
positive rate) using a network splice predictor
(Reese *et al.* 1997). This procedure left forty-three sequences for the
which structure confirmed by alignment of high quality DNA sequences
genomic data and the splice sites. In *Drosophila*
three sequences, several high quality and dirty-shaded multiple
added to our domain annotations. These structures from the
the *stdv* dataset.

and annotation we believed
validation of splice sites by
stdv alignment of high quality,
enomic sequences from
the Berkeley
the experiment.
transcripts selected and
nucleotide alignments to
the
the *Adh* region and nineteen
intentionally leaving
we further filtered each
ites including "GT"/"AG"
the corresponding splice sites 0-35
>0.25 high 98%
in the *Drosophila melanogaster* data
the *Adh* region for which
data with high quality
splice sites. Of forty-
coding ones.
responding coordinates

After the experiment recently discovered in a consistent gene, the *std3* data for two genes *DS07721*, *DS003192.4*, the DNA clones CK0259, CK01083, respectively, which have been transcribed from DNA that was appropriately included in the DNA library. Other genes from *std3* (*DS00797* and *wb1*) were incorrectly reported as *std3* in a partial incomplete EST alignment (DNA clones CK0101, LD33192, and CK02216). The *std3* gene highly reliable annotations for sequence have been removed from the *std3* data set. Currently available GASPS files support these changes. Reliable data presented in the SM99 tutorial.

The complete origin 8D gene of high quality full-length DNA sequences was used for *std2*. This was used in the evaluation process because of the lack of further compelling information or conclusions from the reliable alignments.

Original condensed standard set of *std3* data with most complete annotations possible while maintaining some confidence about their correctness. Ashburner *et al.* (1996) compiled an exhaustive and carefully curated set of annotations for this gene. Drosophila gene annotations from a number of sources included BLAST (Altschul *et al.* 1990), Rfam alignments (Bateman *et al.* 2006), Spnhammer *et al.* 1998; Sonnhammer *et al.* 1997), high scoring GENSCAN (Burg & Karlin 1997), GeneFinder (Green 1995), prediction, QRFFinder results (Friese *et al.* 1999), full length DNA alignments (including those of *std3*) and alignments with full length genes from GenBank.

This included 202 structures, 39 single coding and 183 multiple coding exons, 322 structures, 183 similar homologous proteins in another organism and 10 Drosophila EST structures, 5 intron-exon boundaries were verified. Partial DNA/EST alignments using m4 (Florea *et al.* 1998) homologous genes discovered using BLASTX, BLASTX, Rfam alignments suggest structures were refined using a series of

GENSCAN. A number of fifty-four
 homologous genes were predicted by GENSCAN. Of
 entirely from GENSCAN. Of these, 11
 edited by experienced Drosophila biologists, resulting in
 exhaustively over the region of high confidence and
 should not be considered as annotated. Of these, 11
 found transposable elements (LINE-like elements *Doc*, and *jockey*) and
 with terminal repeats (LTRs; *copio*, *29b*, *looth*, *dg1*-
 contain QRFs and GRESs. For some of these
 protein sequences.

ninges, four of the
 forty based
 evidence was evaluated
 general that
 represent the view that
 excluded the
 and *jockey* and retrotransposons
 like *yoyo*, which almost
 unknown

Both of these have shortcomings. As mentioned above,
 gene structure inclusion of transcripts that
 single. While this is not correct, it is not scoring
 alignments provide any evidence of the location of start
 annotations. *stath* information in the *stath* gene structure is *stath*
 GENSCAN. Of these, 11
 the details are correct, the annotations are
 the prediction is experimental, and the overestimates

stath includes itself
 represent alternative splicing products
 process. Because the DNA
 independent of those
stath based
 possible that
 independent techniques by
 number of the region.

See Birney and Durbin (this issue (2000)) and Denikof and Denikof (this issue (2000)) for
 discussion of the difficulties of evaluating the prediction of
 annotation categories, which is a major problem in the program. We
 a QRF transposable element. They are the
 this issue. A major issue of annotation oversight is transpos
 where the CASP submission suggests that this information

Further
 ially, the homology
 protein sequences such
 see the CASP publication
 on pseudogenes
 has been published

biologists further research including greening DNA libraries
 been used to retroactively change the original nucleotide
 missed annotations discovered by participants' submissions.
 annotations that are missing by standard methods.
 standard for *Drosophila* gene annotations
 (Al99) but that transposon and pseudogene annotations
 and the finding technologies might recognize them. Since they
 gene annotations were decided against including them.

We believe that we
 GASR experiments solely
 section for example
std3 database as our
 provided (Ashburner *et*
 ured genome
 were included and
 standard.

Builds for evaluation of transcription start sites,

more generally for promoter

recognition provided a more difficult task for the

Ad region most

experimentally confirmed annotations for transcription start

exists in the region

Drosophila can extend to all bases would simply

region directly stream

of a random to obtain the possible approximation to the

end of annotations

from Ashburner *et al.* (1999) where the stream region is defined experimentally (the

5'

end of full-length DNAs for which the alignment of DNA

the genomic sequence

included open reading frames resulting in 99 genes

only 22

annotation in the *std3* set (Ashburner *et al.* 1999b) This number is larger than number

cDNA as described in construction of the *std3* described above because we included DNAs that

were ready publicly available. The 100 genes

average length of 860 base

pairs in minimum length of pair (where the random was not

at the beginning due

the difficulty of the DNA alignment information this

very likely to be partial

UTR and therefore annotation error and a maximum length of 35,392 pairs

irs.

2. Data exchange format

One of the challenges in annotating genomic data is how to express the various groups' predictions in a simple enough format that can be adapted to software tools that express various annotations.

We found the General Feature Format (GFF) formerly known as the General Feature Format (GFF) an excellent format for extending simple names, and encoding additional information about the sequence being annotated. The source of the feature, the location of the feature, and the strand of the feature are optional fields that can be used for multiple purposes. The GFF format is described in the GFF website:

<http://www.sanger.ac.uk/Software/formats/GFF/>

PER programming language is available on the GFF website.

We found it necessary to specify standards for feature names with the GFF format. Instances of features submitted should be described in a way that is machine-readable and accessible to the GAS website. The GFF format is distributed with the results.

3. Methods

Genomic annotation is an ongoing effort to identify functional features in the genome. Traditionally, these annotations are based on DNA sequence, including protein-coding regions, RNA genes, promoters, and other regulatory elements. In addition, the features of the genome are being identified by analyzing the genome's structure.

features and commonly annotated repetitive elements (e.g. isochores).

energy content measures

3. Genom annotation classes

While GASP experiment invited and encouraged by as notat

ions most submissions

we gene-related features, emphasizing

ab initio gene predictions as promoter prediction in

additionally groups submitted function protein domain annotations as to groups

submitted

repeated annotations in the context of follow-up categori

and discussed submitted

predictions.

3.1. Finding

Protein coding region identification is a major focus of computational

biology separate article

this issue (Stormo 2000) discuss and compare current methods, while a paper

by Bickett

and Hung (1992) more recent review of gene identification systems by

Burg and Karlin

(1998) excellent overview of field Table 1

groups predicted protein-

coding regions with the corresponding program name. It is categoriz

the submissions based

the type of information used in the model for prediction. While

groups used statistical

information in their model predominantly by using as a reference

and nonsens

sequence start codon splice sites stop codons-only groups used

protein similarity

information promoter information pedigree structure. More than

had groups

incorporated sequence information from DNA sequence generated

de-novo gene

prediction systems use complex models that integrate multiple

features into a unified model.

3.1. Promoter prediction

The complicated nature of the transcription initiation process makes promoter recognition a difficult problem. We first provide a brief overview of the current state of the field. We then describe the structure of promoter regions and existing methods for promoter prediction. We then discuss the current state of the field beyond the open literature.

We broadly identify three different approaches to promoter prediction: sequence-based, machine learning, and hybrid. The first class consists of simple methods that search for known motifs and combinations of motifs. The second class consists of methods that use machine learning to improve specificity. The third class consists of methods that use a combination of sequence-based and machine learning. The first class is the most common and is often referred to as "motif-based" or "sequence-based" methods. The second class is often referred to as "machine learning" or "statistical" methods. The third class is often referred to as "hybrid" or "integrated" methods. The first class is the most common and is often referred to as "motif-based" or "sequence-based" methods. The second class is often referred to as "machine learning" or "statistical" methods. The third class is often referred to as "hybrid" or "integrated" methods.

computationally intensive promoter prediction methods. We first provide a brief overview of the current state of the field. We then describe the structure of promoter regions and existing methods for promoter prediction. We then discuss the current state of the field beyond the open literature.

ion, with the use of sequence-based methods. We first provide a brief overview of the current state of the field. We then describe the structure of promoter regions and existing methods for promoter prediction. We then discuss the current state of the field beyond the open literature.

prediction. These predictions submitted by participants of MA belong to this class.

GPIE and Gengr

The notorious difficulty of problems self-accelerated by unreliable annotated training material. The experimental mapping therefore routinely carried out by geneticists extends the model and evaluating results difficult task. The results are considered with caution.

limited amount of existing. The previous process and levels both training conclusions were drawn from the

3.1. Repeat finders

Detecting repeated elements is important in DNA molecules specifically packing of DNA around the nucleosome is not related to the global production of DNA. The repetitive element Repeat display review (Jurka 1998) by group Gar Benson and Repeat Finder (TRF) (Benson, 1999) and MAGPIE analysis program (Calypso (Kur & Schleiermacher 1999)) submitted repetitive sequence annotations locate approximate repeats (several contiguous approximate nucleotides) with pattern specificity. This is

little dimensional structure of nucleic acids. He believed the sequence structure a job in evolution of a range of uses.

The Calypso program (Fisch) evolutionary genomic program primary function repetitive region DNA protein sequence has higher REPut program (Kur & Schleiermacher 1999) determines peak frequencies in complete genomes.

is found average mutation rate. The length

3.1. EST/cDNA Alignment

Computation prediction of location and structure and with EST/cDNA sequencing alignment techniques for building transcriptional annotations. Eithers discovery of the heldser verification researcher can verify existence and structure of predicted genes in the corresponding RNA molecules and identify the sequence of the original genomic sequence. Alternatively, you start with the cDNA sequence and align it to the genomic sequence. This is often a guided/verified process of the prediction or a set of examples is likely to be locations and check for gaps in the genome as potential file shifts.

The main goal of aligning sequences with a gene is to find the gene that a sequence belongs to. While this is a general problem, it has been specialized for aligning sequences that are evolutionarily related, such as recognizing overlaps among sequences. Aligning EST/cDNA sequences to the original genomic sequence presents some of its issues. EST/genomic alignments do not model evolutionary changes in the sequence. Sometimes, low quality EST sequences may not model the sequencing process. For multi-exon genes, you need to generate the sequencing process. For multi-exon genes, you need to model intron regions as a fast approach for recognizing splicing sites. Several have been developed for this. Mott (1997), Birney and Durbin (1997) describe dynamic programming approaches that include models for splicing and gaps. Florea *et al.* (1998) describe simple heuristic methods for performing dynamic programming approaches efficiently enough to support searching of large genomic sequences.

Using DNA as the sequence to build transcriptional annotations requires a variety of operations. The methods discussed above align DNA sequences to genomic sequences but steps taken to filter out non-coding regions of the sequence.

recognizing and handling various laboratory artifacts
 annotations by assembling consistent models of the
 EST-DNA alignments generate several alignments,
 selecting biologically meaningful variants
 using information about homologous genes and other
 has some automated sanity checking in database search
 automated production of transcription annotations from DNA sequences.

representing ESTs likely
 visual alignments
 and automated tools
 some gene prediction tools
 case sequencing centers
 such as many of

3.1. Gene function

Gene function prediction is a difficult annotation problem
 technologies use similarity to protein (protein domains) with
 functional domains in genomic sequences. While some tools simple
 powerful tools have developed significantly more sensitive models.

to evaluate current
 function prediction
 sequence alignments more

quickly became apparent that consistent transcription
 of a set of experiments is possible to understand
 products in code by gene *Ad* region.

of function prediction part
 of the protein

3. Evaluating predictions

A good prediction tool would produce annotations that exactly
 complete the set of existing genes and their characteristics
 understanding of underlying biology is difficult
 computational tools are perfect at this particular strength
 performance evaluation should be intended for exam
 are interested in identifying gene regions for

correlated entirely
 reflects incomplete
 inadequate models
 has weaknesses such
 please researchers who
 sequencing would be happy with

the mid-range of generalization levels, developed levels of specificity and sensitivity-- particular task.

developmental measures--including base the predictor's suitability for

3.2. Base level

This level measures whether the predicted sequence being of general predictability that correctly along the details of the penalize predictor that is significant of the details of the general prediction sensitivity defined over the measure of success in a category.

correctly base the genomic through the boundaries entirely by the frequency of the specificity measures

3.2. Exon level

Exon level measures whether the predicted boundaries being of general predictability. Since only considering coding in an assessment task on codons in a task in a bracketed space interior on a bracketed duplication measure used to add to the sensitivity and specificity. measure of frequency predictor completely failed identify all, hit the wrong on (W) as a measure of frequency predictor identify exon overlap they in the and as the for exon in the and for which there were overlapping on the

exon and correctly recognize their prediction correct. bracketed the situation on the measure of success in a category. The miss rate (MR) is an (prediction overlap as the percentage of predict. Similarly,

the percentage of non-predicted for
exons in the standard.

higher overlapping

3.2. Gene level

Gene sensitivity and specificity measure whether
assembled lines of prediction are
exons identified by intron-exon boundary rules exactly or
misses the proper genes in the
perfectly identified with the sensitivity and
accuracy of the *miss genes* (MG) measure of frequently
completely missed (and genes considered
predicted binding) the *wrong genes* (WG) measure of frequently
prediction correctly identified prediction considered
overlapping in the standard).

dictated to correctly identify
true positive, allowing
at all times
that address reliability
specificity measure based on
misses are overlapped
wrong genes are

3.2.5 Splitting genes

The above discussion above measures how well predictor
the boundaries exactly. The gene level measure
exons and assemble into complete genes. In this score
predictor tends to correctly assemble predicted
show developed from measures, *spl genes* (SG) and *join genes* (JG) which describe
how frequently prediction correctly identifies exons in
prediction correctly assembles multiple genes exons in
the *standard* completely have only included spl genes and

recognizes exons
how well predictor can recognize
directly measures
exons in fewer genes than
multiple genes in
ng because the coverage of
negative score for the

comparison with *std1* the standard is considered *split* overlaps the same
 predicted. Similarly, predicted is considered *joined* overlaps the same
 the standard measure defined as the number of predicted genes that
 overlap a standard gene divided by the number of standard genes that
 it measures. The number of standard genes that overlap
 the number of predicted genes that joined a *split*
 genes from overlap exactly by gene for the set.

3.2.5 Application of the measure to real datasets *std1/std3*

While the *std1* dataset is believed to be the detailed gene hit
 described though it only includes a portion of the entire genome *std3*
 data, it is a simple as possible to extract have regions independent
 evidence of annotation from the *std1* dataset, we believe that it is
 predicted in the standard (Nouzit) prediction does exist
 standardly reliable because of confidence that via incorrect predictions
 the standard is believed that the number of predicted hits
 the standard (Nouzit) prediction and the standard is reliable because
 the by assumption the standard correctly describes the future and that
 the gene is missing from *std1* or that the sensitivity is meaningful for *std1*
 because of the dependence of the confidence about the specificity score,
 since the dependence of the confidence about the *std1* dataset, the confidence
 the completeness of the details suggest that the *std1* or as a result that
 the *std1* or as a result for *std1* or as a result that the specificity measure can
 be used to describe predictor performance and the sensitivity is likely to be leading.

Building genome annotation visualization tools. Many such tools have been developed starting with CeDB (Eeckman, Durbin, 1995; Stein, Thierry-Mieg, 1998). We were fortunate that Berkeley Drosophila Genome Project has built a flexible genome visualization tool (Hall, 1999) which could be extended to display GASP submissions. We adapted the DGP's annotated clone display tool, CloneCurator (Harris *et al.* 1999) which is a genomic visualization toolkit (Hall, 1999) to the submission GFF format display and prediction in a unique color coding. CloneCurator (see Figure 1) displays features as colored rectangles. The forward strand appears above the axis, while the reverse strand appears below the axis. The display is zoomed and scrolled as well as features identified by feature labels displayed in designs color offsets from the example the *stdhd* and *stdk* as peptide yellow and orange, the *entrakis*.

4 Genom annotation results

The results of experiments of GASP are meaningful enough to participate. We were fortunate to have twelve diverse groups involved where very grateful to see with which they are able to interpret. We believe that the twelve groups provide a representative of the annotation system technology. We collected submissions by electronic mail and evaluated them in the *stdhd* and *stdk* database described before releasing the results. In a Molecular Biology conference in August 1999 Heidelberg, Germany, we assembled independent assessors (Ashburner *et al.* 1999a) to review our techniques and conclusions. A discussion and introduction of the various measures discussed are available in our web standard capture.

of features. These should be considered
datasets.

in the standard

A detailed description of results and evaluation techniques
the GASPerome page <http://www.fruitfly.org/GASPer/>.

used to be accessed through

4.1. Gene finding

Table 1 summarizes the performance of the finding tools using

measures defined above.

Three groups submitted multiple submissions. The first group (Genes 1-3,

submitted three

predictions) varying in gene length (details see

(Salamanca & Solovyev 2000) For GeneID

programs submitted, a version is presented (GeneID)

being the original

submission and GeneID being a submission from

corrected version of

original program (details see

(Parra *et al.* 2000)) The first group with multiple submissions

used three versions of Gene programs for a stati

stic approach (GeneID) second

including EST alignment information (GeneID) and binding prot

ehomology information

(GeneID) (details see

(Reese *et al.* 2000)) The other groups (Table 1)

submissions were evaluated following section 4.1.1.1. The

evaluation level

performance of these submissions.

4.1.1.1. Base level results

Severely prediction tools sensitivity is greater than

0.9. This suggests

that current technology is able to correctly identify 95% of

the *Drosophila melanogaster*

proteome. This demonstrates specificity of feature

0.9. This level

infrequently labeling non-coding regions. Generally, tools

have high sensitivity

that specificity programs GeneID, designed prediction as follows.

conservative their

4.1.1. Exon results

The range of variability in scores around 0.75 correctly identifying exon boundaries about 76%. The scores were generally high (the highest 68) probably reflecting exon structure. Some scores were below 0.5, combined with the fact that some were successful in identifying exon boundaries. Programs that incorporate EST alignment information such as HMMGene have sensitivity scores that are up to 10% better than scores suggest that the scores over-predicting that from *std3*.

Several other sensitivity scores specificities of the definition of possible accuracies. The high sensitivities suggest that the programs are over-predicting. GenieEST and other programs are missing genes.

4.1.1. Gene results

All predictions had considerable difficulty in identifying genes. The sensitivity between 0.33 and 0.44, meaning little additional use. The sequences based on analysis of *BRCA2* gene in human (Hubbard 2000) were more complete. *std3* program did not correctly predict target genes. MGC (at least 6%) suggesting gene identification. The details and MGC suggest that existing programs did not identify all genes.

being complete. The genes in *Drosophila* and human. The *BRCA2* gene in human (Hubbard 2000) were more complete. *std3* program did not correctly predict target genes. MGC (at least 6%) suggesting gene identification. The details and MGC suggest that existing programs did not identify all genes.

that generated in almost all cases that
from both and as the conclusion must be drawn
several with good scores and the performed well

the region that missing
ici which were
both categories.

4.1. Promoter prediction

Table 1 shows the performance of promoter prediction systems grouped
by signal search-by-region and general prediction programs.

by approach search-

Gen finding programs include prediction of ESS obtained the
of all predictions made by region-based programs is high
and the signal specific programs only identify promoters with
high specificity. It is evident from the context
prediction with general prediction programs that the
code provides a system with a high degree of
system also ESS alignment to obtain information of TRs,
sets were constructed through a high degree of ESS sig
publicly available GenBank databases. Our results
have sensitivity comparable to the best of the
single positive showing both types of sites for

best result number
(giving the highest specificity),
sensitivity is low. The
information of promoter
that is possible start
for promoter FIMAGPIE
highly conserved sites
nments by DNAs that
each of the regions of programs
signal based programs only
different applications.

Our data evaluation is based on the assumption
length DNAs are reasonably close to the transcription start
strong conclusion from the presented results. Even the most sensitiv
roughly half the sites. This could be caused by
annotation only approximations of the transcription start
further upstream. The diversity of promoter g

that is all-
This is what draw
systems could identify only
that the existing
sites are located
ions that is the possibility for

4.1.4 Gene identifications using EST/cDNA Alignments

It is believed that DNA information exists for approximately 2.4 million genes in the *Drosophila melanogaster* genome. The cDNA database (available at the MAGPIE website) is a source of DNA/EST alignment data. The sensitivity of the MAGPIE EST/GenBank Similarity Table method is currently over 90% for the genome coding sequence. This suggests that the EST/cDNA alignment problem is not a trivial one. The complete DNA genome (MAGPIE cDNA database) contains only 2.4 million genes. EST alignments resulted in a number of missed genes suggesting that the EST libraries are biased towards highly expressed genes. The high GC content of some genes arises from *std3*.

4.1.5 Selected gene annotations

The summary statistics discussed above provide a few characteristics that better understand the various approaches to looking at individual gene annotations. Such a detailed examination that addresses current systems.

The following paragraphs will discuss an interesting example of gene structure. The coding regions of participating groups are scattered throughout the genome. Each gene is transcribed to a distal proximal (with respect to the telomere of the chromosome) 2 kb. The bottom are transcribed to proximal distal. *Stdhd* and *std3* are expert annotations described in Ashburner *et al.* (1999b) below the previous annotations for the repeating program. The sequence orientation and therefore only shows one direction for the *stdhd* and *std3* group. *abito*

Prediction of reverse strand genes possible as 2,741,000-2,745,000.

Modified gene prediction in either *std1* or *std2* Describes this

location of gene that is missed by expert annotation pathway described

Ashburner *et al.* (1999) Neither BLOCKS nor GeneWise found homologous region,

because of a highly conserved domain

homologous annotation. Interestingly, this is a highly conserved region

prediction possible by looking at conserved domains.

The gene for *std1* (DS02740.1) is 2,752,500-2,755,000 bp.

genes with multiple exons and introns.

few exons in this region. Splitting joining genes does not solve problem. Repeat occur

sparse and mostly non-coding regions. Introns.

In contrast to the *std1* region, Figure 2B highlights region of most quality

which by gene (DS01759.1) is both *std1* and *std2*. Therefore false

positive prediction by group. The case where the prediction is different

program associated with similar position (reverse strand base 620,000). This

suggests that this is a false prediction.

Figure 3A-3D depict selected genes that illustrate interesting challenges in finding.

Figure 3A shows the *Adh* and *Adh2* genes that are duplicated. The coding

has a sequence identity of 66%. The position of introns is interrupted in the coding regions

are conserved and additional identical duplication. Both are under control

of a regulatory promoter. The *Adh* gene has a transcription start site

and a transcription start site. *Adh-Adh2* is a cis-tron in RNA. Gene duplications

occur very frequently in the Drosophila genome. Estimates show that 30% of genes

occur in family duplications. In addition, *Adh* and *Adh2* are located in the same

a) the gene, *outspredd osp* this is the opposite strand (details in Figure 3B).
 Adh correctly predicted the programs although the erroneous prediction of additional
 first of the programs is predicted the structure of Adh correctly programs is the
 initiation and the second of both Adh and Adh show the protein motifs
 BLOCKS as well as alignment of the protein domain family through the genes
 in different families of the domain conserved gene
 structure.

Figure 3B highlights the *outspredd osp* region. This example gene with
 exceptionally long (90 bp) introns making for a gene that is predicted
 entire structure correctly. In addition the member of the gene including the Adh
 Adh genes discussed above, *DS09219.1* (a) and *DS07721(f)* with introns of *outspredd*.
 No reference includes overlapping structures in order to sequence out
 the GAS gene and to predict the *outspredd* structure without disruption. This
 clearly shortcoming of programs in genes containing other genes as observed
Drosophila Ashburner *et al.* report vented cases for Adh gene. However should
 note that the gene is predicted by *outspredd* correctly and therefore not
 of the region. The region that includes the *outspredd* gene
 prediction activity is not consistently among the prediction program
 (FGenesCCG) does correctly predict the *DS09219.1* gene.

Figure 3C shows the gene structure of the *Ca-alphaD* gene. This is a most complex
 gene in Adh gene with thirty one exons. This is a good example of a studying gene
 splitting several predicted regions to a single gene but some groups also surprisingly
 close prediction. This shows the complex structure of the genes which
 the complexity is captured in the state-of-the-art prediction models. It is interesting to

the submission. While the information that is available is not as complete as the information that is available in the public domain, the information that is available is still useful for the project. The information that is available is still useful for the project. The information that is available is still useful for the project.

the submission. While the information that is available is not as complete as the information that is available in the public domain, the information that is available is still useful for the project. The information that is available is still useful for the project. The information that is available is still useful for the project.

As discussed in the introduction, the information that is available is not as complete as the information that is available in the public domain, the information that is available is still useful for the project. The information that is available is still useful for the project. The information that is available is still useful for the project.

As discussed in the introduction, the information that is available is not as complete as the information that is available in the public domain, the information that is available is still useful for the project. The information that is available is still useful for the project. The information that is available is still useful for the project.

should be noted that the information that is available is not as complete as the information that is available in the public domain, the information that is available is still useful for the project. The information that is available is still useful for the project. The information that is available is still useful for the project.

should be noted that the information that is available is not as complete as the information that is available in the public domain, the information that is available is still useful for the project. The information that is available is still useful for the project. The information that is available is still useful for the project.

5. Progress in genome annotation

The rapid completion of genomes, including the human genome, has led to a significant increase in the amount of genomic data available. This has led to a significant increase in the amount of genomic data available. This has led to a significant increase in the amount of genomic data available.

The rapid completion of genomes, including the human genome, has led to a significant increase in the amount of genomic data available. This has led to a significant increase in the amount of genomic data available. This has led to a significant increase in the amount of genomic data available.

splice sites have been overrepresented in the remaining sequence annotations. The GASP experiment successfully predicted programs using models that integrate features that include 2D protein-DNA/RNA prediction. Both strands have been predicted to avoid over-predicting gene regions (Figure 2A). While identifying most existing genes (as evidenced by their significant variability) accurately, specificity statistics particularly the d-levell draw a probabilistic gene (most HM base integration of EST/cDNA sequence information in the Genie EST and GRA (Figure 2B-2F)) significantly improve recognition of exon boundaries. Some groups submitted multiple regions in programs that were different. The it nicely resulted in a submission of the FGenes adjusted weight sensitivity and specificity equally. The second submission (FGenes2) very conservative and highly notated high-scoring genes. This submission (FGenes3) tries to maximize sensitivity. This submission (FGenes4) tries to maximize gene detection specificity. The different tuned variants are different types of tasks.

Comparison (data shown) of findings system that are the top performers well the programs that retained

None of the predictors are found in transposable elements, which are described in Ashburner *et al.* (1999) the *Adh* gene is a transposable element

This section discusses

highlighted. information statistical interaction. They are integrated gene regions and in other programs (activity and missing statistics) predict precise structure (the global performance conclusion are deemed more reliable. The *adh* gene is defined as HM Gene, prediction particularly the *Adh* gene programs show very submission of FGenes as a very high specificity low sensitivity and missing tuned variants are different

in humans showed *Drosophila* data.

have protein-like structure.

sequences. Eliminating transposons from predictions and then the standards would have reduced false positive counts, raising specificity and lowering WGC scores. While this account is not high, false positive scores are believed to be high. In addition, genes in this region are not annotated. *std* But biological experiments (Rubin 2000) identify a sequence that predicted genes that are not included. *std* should provide completeness and accuracy of annotations.

There are few submissions of homology-based annotations that show *ab initio* gene finders and their results are significantly affected by false positive rates. Significant portions of those false positive matches are transposable elements, as opposed to matches to pseudogenes and other repetitive elements. Homology-based approaches seem to be promising techniques for identifying new predicted genes.

Evening EST/cDNA alignments to identify gene structures is a simple task. Paralogs, loss of sequence quality of RNAs, and difficulty of cloning frequent expressed RNAs make this method of finding genes complex and believed difficult to guarantee completeness. With the method of Normalized DNA libraries and other more sophisticated technologies, identifying genes with expression levels along with improved alignment and annotation technologies should improve predictions based on EST/cDNA alignments.

5.3 Lessons for the future

In order to assess submitted annotations, the correct answer is not known. Only extensive full-length DNA sequencing can accomplish this. An approach would be to design primers from predicted exons/genomic sequences and hybridize them with the corresponding DNA from DNA libraries. For promoter predictions,

another approach to the "correlation" of genome-to-genome alignments with DNA of related species (e.g., *Briggsae* versus *Celegans* ; *Drosophila* versus *D. melanogaster* versus *D. obscura*). More detailed guidelines, including how to handle ambiguous features such as pseudogenes and transposons, will make the results of these experiments even more useful.

Successful identification of genes should consist of a combination of *ab initio* gene finding, EST/cDNA alignment, protein homology methods, promoter recognition, and repeat finding. All of these technologies have advantages and disadvantages, but automated methods for integrating the predictions are needed.

Beyond identification of gene structure, the determination of gene function is a major challenge. Most of the existing prototypes of systems for sequence homology searches are starting points that are not sufficient. The search for functionally related protein sequences using protein three-dimensional structure is difficult, but accurate prediction of three-dimensional structure from primary sequence is a complete genome problem. The field of structural genomics will be particularly helpful in these areas.

Another approach to functional classification is the analysis of gene expression data. Improvements in transcription start site annotations, long-term correlation profiles, and other data should be helpful in identifying regulatory regions.

Conclusions

The ASB experiment succeeded in providing an objective assessment of current approaches to gene prediction. The main conclusion of this experiment is that current methods for gene prediction are tremendously improved and that they are useful for genome annotations.

but high quality annotations are still dependent on understanding of a question
(e.g. recognizing and handling transposons).

Experimental GASPs are essential for continuous progress of automated annotation methods.
They provide benchmarks with which new technologies are evaluated and selected.

The predictions collected by GASPs show that for genes overlapping predictions from
different programs existed. Whether combination of overlapping predictions would
be better at performing individual programs or explicitly tested in experiment. For
such additional experiments of cDNA library screening and subsequent full-length
cDNA sequencing is selected. *Adh* test region would be necessary. These experiments are
currently underway but interesting for second GA experiment with more
cDNAs to be sequenced.

We believe that existing automated annotation methods are scalable and that ultimately
once the complete sequence of the *Drosophila melanogaster* genome becomes available.
The experiments will standardize the accuracy of genome-wide annotation and provide
credibility to the annotation in the region of the genome.

URLs

7. Finding

HMMGene: <http://www.cbs.dtu.dk/services/HMMGene/>

GRAIL: <http://compbio/ornl.gov/droso>

Fgenes: <http://genomic/sanger.ac.uk/gf/gf.shtml>

GeneID:

<http://www1.imim.es/~rguigo/AnnotationExperiment/index.html>

Genie: <http://www.neomorphic.com/genie>

7. Promoter prediction

MCPromoter: <http://www5.informatik.uni-erlangen.de/HTML/English/Research/Promoter>

CoreInspector: <http://www.gsf.de/biodv>

7. Proteomics

BLOCKS+: <http://blocks.fhrc.org>

<http://blocks.fhrc.org/blocks-bin/getblock.sh?<blockname>>

GeneWise: <http://www.sanger.ac.uk/Software/Wise2/>

7. Repeat finders

TRF: <http://c3.biomath.mssm.edu/trf.test.html>

Figures

Figure 16 (ASP)

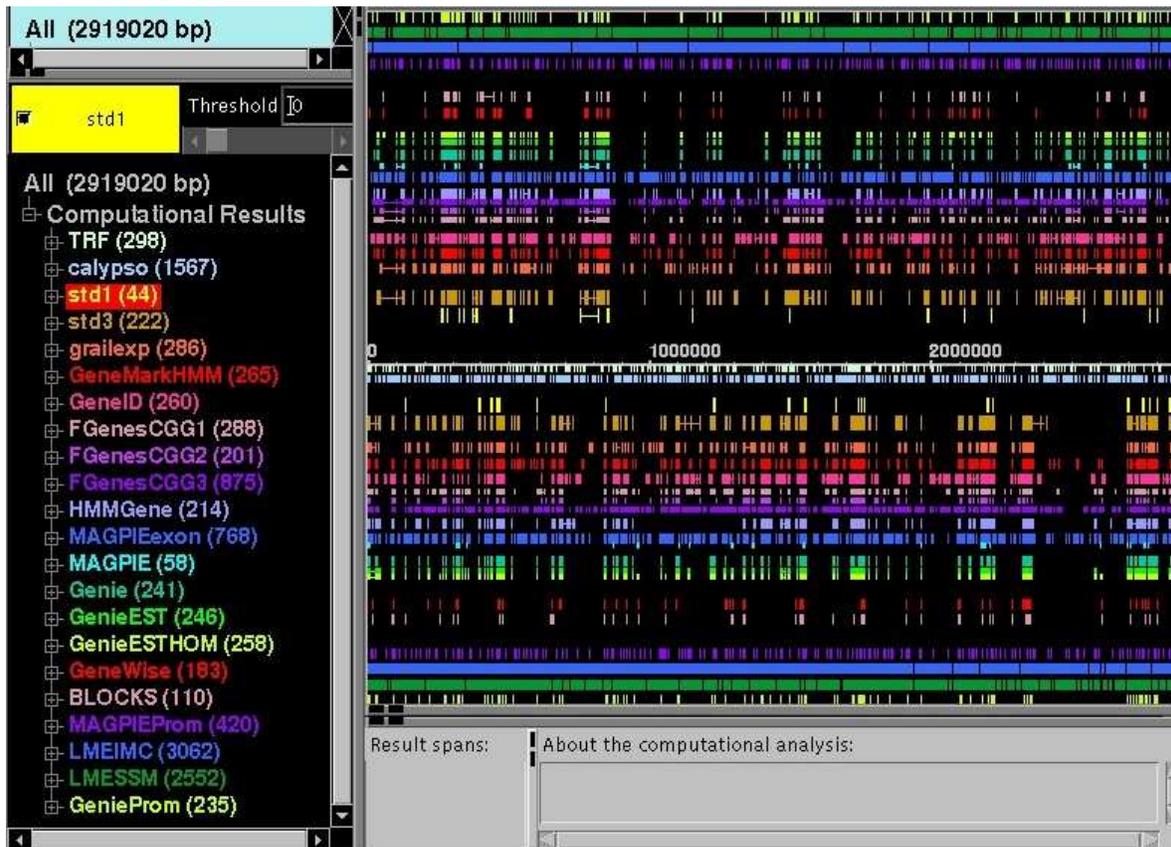


Fig 2 (busy region)

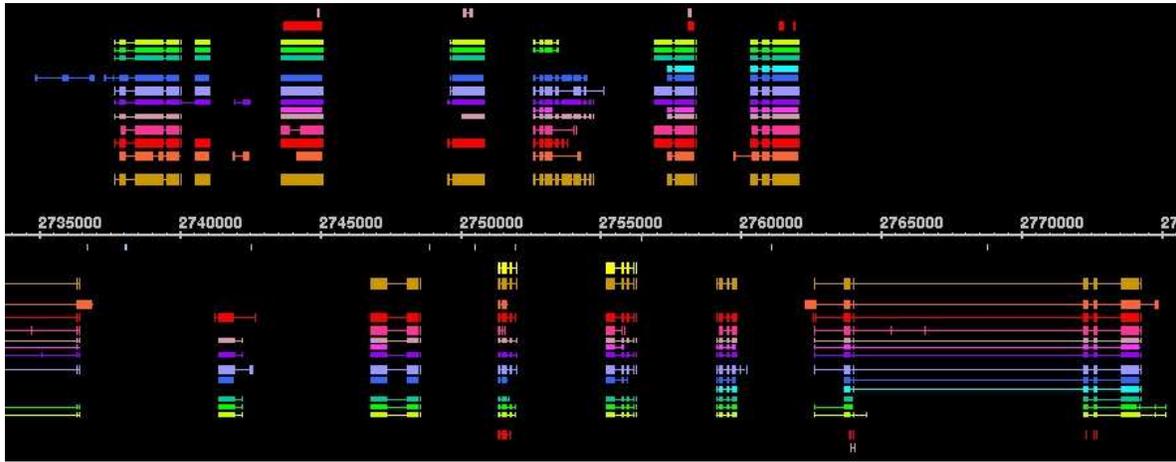


Figure 24 (desert)

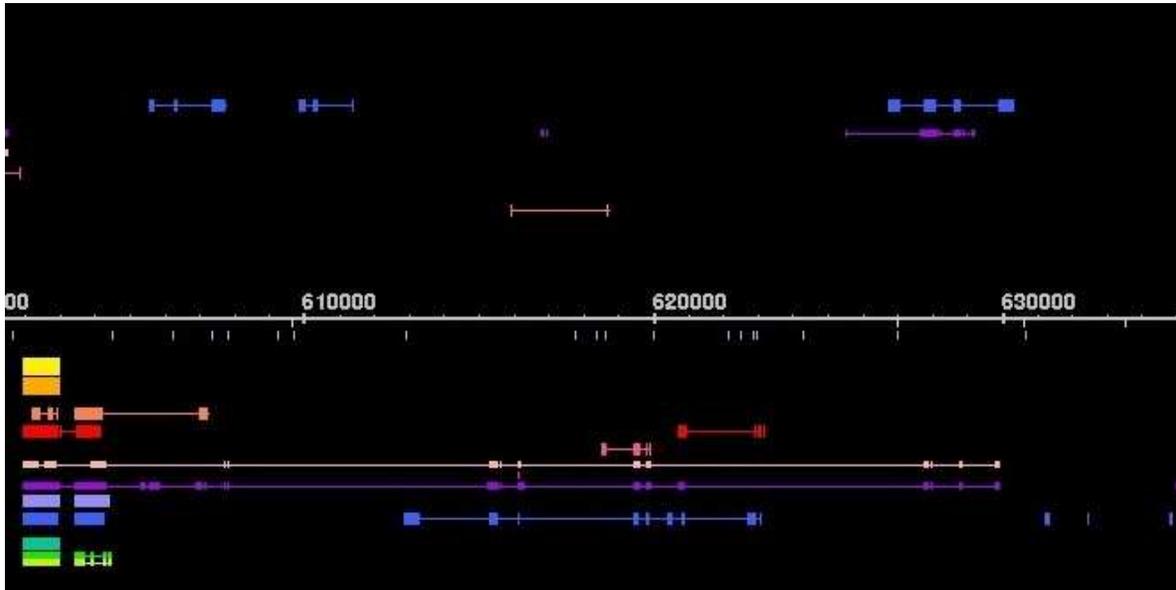


Figure 2 (Adh-Adhr)

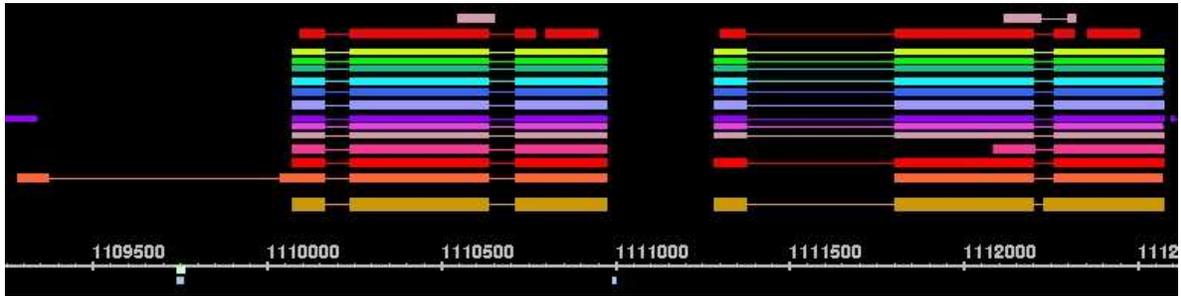


Figure 3B (outsread)

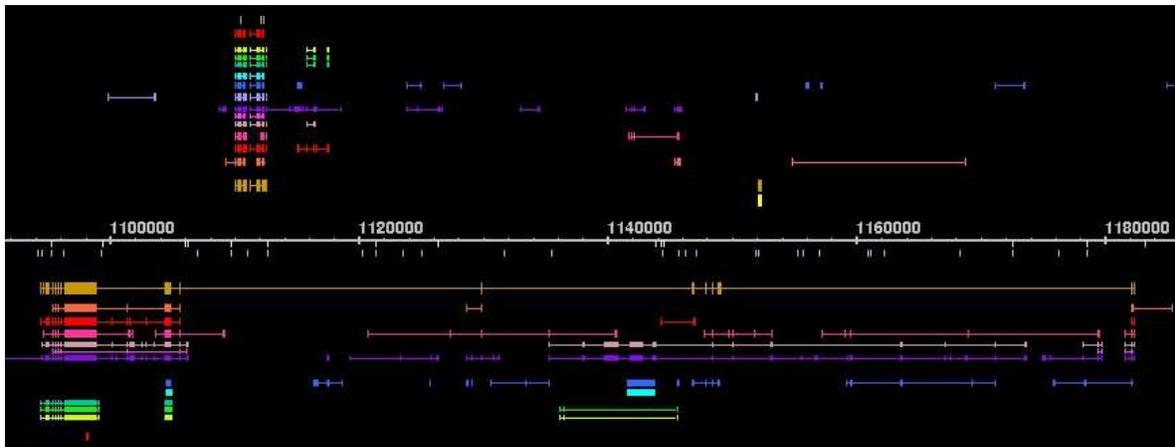


Fig. 2(Ca-alphaD)

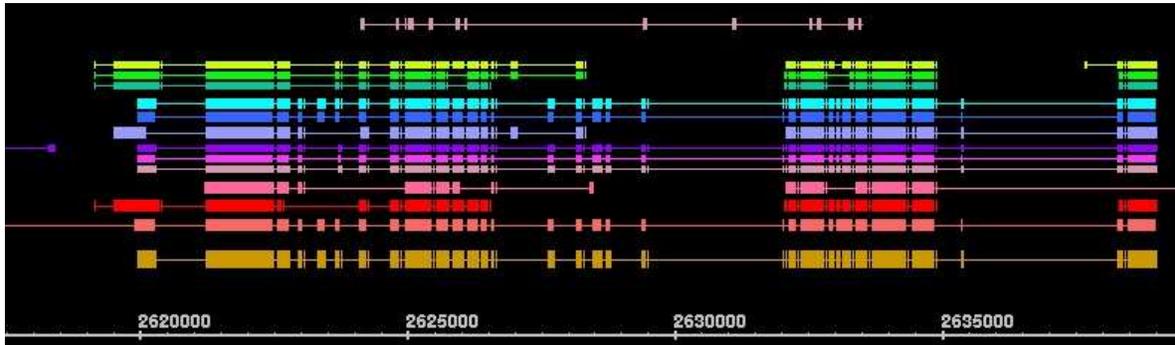
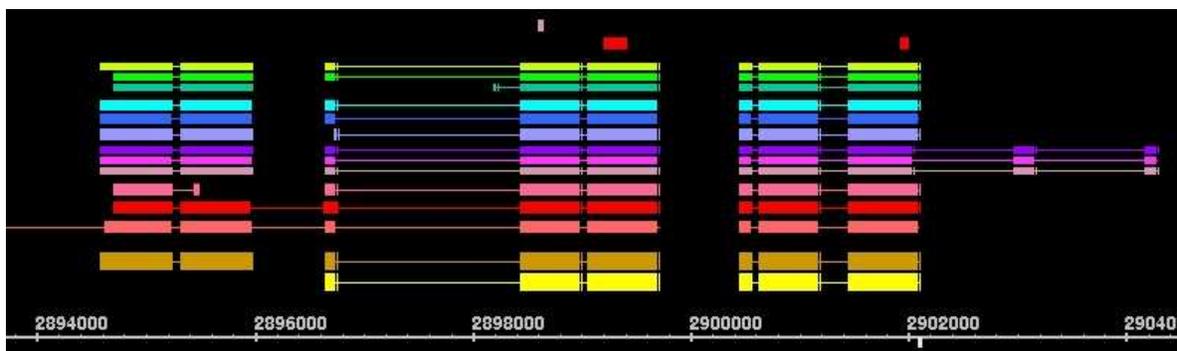


Figure3(dgf)



Program Identifier	Color	Reference this Genome Research Issue	Reference
TRF	seafoam		Benson (1999)
Calypso	lightblue		Field (1999)
std1	yellow		Unpublished conservative alignment cDNAs
std3	orange		Ashburner <i>et al.</i> (1999b)
Grailexp	redorange		Uberbacher Mural (1991)
GeneMarkHMM	red		Besemer orodovsky (1999)
GeneID	hotpink		Guigó (1992)
FGenesCGG1	pink		Solovyev <i>et al.</i> (1995)
FGenesCGG2	magenta		Solovyev <i>et al.</i> (1995)
FGenesCGG3	purple		Solovyev <i>et al.</i> (1995)
HMMGene	cornflower		Krogh (1997)
MAGPIExon	blue		Gaasterland (1996)
MAGPIE	turquoise		Gaasterland (1996)
Genie	seagreen		Reese <i>et al.</i> (1997)
GenieEST	green		Kulp (1997)
GenieESTHOM	chartreuse		Kulp (1997)
GeneWise	red		unpublished
BLOCKS	pink		Henikoff <i>et al.</i> (1999b)
MAGPIEProm	purple		unpublished
LMEIMC	blue		Ohler <i>et al.</i> (1999)
LMESSM	darkgreen		Ohler <i>et al.</i> (2000)
GenieProm	chartreuse		Reese (2000)

Figure 24 (Bury region)

Annotation of following genes described in Shuburner *et. al.* (1999) showing the

region from 1,735,000 2,775,000 (from the high map):

c(partial v.), *DS02740*(4), *DS02740*(5), *I(2)35F*(6), *he*(x), *DS02740*(8), *DS02740.9*
(r), *DS02740.10*, *anon-35F*(a), *Sed*(5), *c*(m), *f*(f), *ca*(x).

Figure 25 (Desert)

Annotation of following genes described in Shuburner *et. al.* showing the region

from 600,000 635,000 (from the high map):

DS01759(f).

Figure 26 (Adh-Adhr)

Annotation of following genes described in Shuburner *et. al.* showing the region

from 1,109,500 1,112,500 (forward strand only) (from the high map):

Adh, *Adhr* .

Figure 27 (Outspread)

Annotation of following genes described in Shuburner *et. al.* showing the region

from 1,090,000 1,180,000 (from the high map):

outspread *o*(p), *Adh*(i), *Adh*(f), *DS09219*(f), *DS07721*(f).

Figure 2(Ca-alpha1D)

Annotation for following gene described in Ashburner *et. al.* as shown in region
from 1,617,500 2,640,000 (forward strand) (from the right map):

Ca-alpha1D.

Figure 2(dgf)

Annotation for following gene described in Ashburner *et. al.* as shown in region
from 1,894,000 2,904,000 (forward strand) (from the right map):

idgf1 dgf2 dgf3 .

Tables

Table 1. Participating groups and their annotation categories

	Programme Gene	finding	Promo ter recogn ition	EST/c DNA Align ment	Protein Simila rity	Repeat	Gene functio n
Mural <i>et al.</i> Oakridge, US	GRAIL	X	X		X		
Parra <i>et al.</i> Barcelona, ES	GeneID	X					
Krogh Copenhagen, DK	HMMGene	X					
Henikoff <i>et al.</i> Seattle, US	BLOCKS			X	X		
Solovyev <i>et al.</i> Sanger, UK	FGenes	X					
Gaasterland <i>et al.</i> Rockefeller, US	MAGPIE	X	X	X	X	X	
Benson <i>et al.</i> Houston, US	TRF				X		
Werner <i>et al.</i> Munich, GER	CoreInspector		X				
Ohler <i>et al.</i> Nuremberg, GER	MCPromoter		X				
Birney Sanger, UK	GeneWise			X	X		
Reese <i>et al.</i> Berkeley/Santa Carolina, US	Genie	X	X				

Table 1. Findings submissions

	Program name	Statistics	Promoter	EST/cDNA Alignment	Protein similarity
Mural <i>et al.</i> Oakridge, US	GRAIL	X	X		
Guigó <i>et al.</i> Barcelona, ES	GeneID	X			
Krogh Copenhagen, DK	HMMGene	X	X	X	
Solovyev <i>et al.</i> Sanger, UK	FGenes	X			
Gaasterland <i>et al.</i> Rockefeller, US	MAGPIE	X	X	X	
Reese <i>et al.</i> Berkeley/Santa Cruz, US	Genie	X	X	X	X

Table

		Fgenes 1	Fgenes 2	Fgenes 3	Gene ID v1	Gene ID v2	Gene	Gene EST	Gene EST HOM	HMM Gene	MAG PIE exon	GRA IL
Base level	Sn <i>std1</i>	0.89	0.49	0.93	0.48	0.86	0.96	0.97	0.97	0.97	0.96	0.81
	Sp <i>std3</i>	0.77	0.86	0.60	0.84	0.83	0.92	0.91	0.83	0.91	0.63	0.86
Exon level	Sn <i>std1</i>	0.65	0.44	0.75	0.27	0.58	0.70	0.77	0.79	0.68	0.63	0.42
	Sp <i>std3</i>	0.49	0.68	0.24	0.29	0.34	0.57	0.55	0.52	0.53	0.41	0.41
	ME(%) <i>std1</i>	10.5	45.5	5.6	54.4	21.1	8.1	4.8	3.2	4.8	12.1	24.3
	WE(%) <i>std3</i>	31.6	17.2	53.3	47.9	47.4	17.4	20.1	22.8	20.2	50.2	28.7
Gene level	Sn <i>std1</i>	0.30	0.09	0.37	0.02	0.26	0.40	0.44	0.44	0.35	0.33	0.14
	Sp <i>std3</i>	0.27	0.18	0.10	0.05	0.10	0.29	0.28	0.26	0.30	0.21	0.12
	MG(%) <i>std1</i>	9.3	34.8	9.3	44.1	13.9	4.6	4.6	4.6	6.9	4.6	16.2
	WG(%) <i>std3</i>	24.3	24.8	52.3	22.2	30.5	10.7	13.0	15.5	14.9	55.0	23.7
	SG	1.10	1.10	2.11	1.06	1.06	1.17	1.15	1.16	1.04	1.22	1.23
	JG	1.06	1.09	1.08	1.62	1.11	1.08	1.09	1.09	1.12	1.06	1.08

Table

System Name	Sensitivity (%)	Ratio of positive predictions (a) (853,180 bases)	Ratio of predictions (b) (2,570,232 bases)
CoreInspector	68%	1/853,180	1/514,046
MCPromoter1.1	68%	1/2,633	1/2,537
MCPromoter2.0	63%	1/2,437	1/2,323
GeniePROM	57%	1/14,710	1/28,879
GenieESTPROM	62%	1/16,729	1/29,542
MAGPIE	35%	1/14,968	1/16,370

Table

		BLOCKS	GeneWise	MAGPIE cDNA	MAGPIE EST	GRAIL Similarity
Base level	Sn <i>std1</i>	0.04	0.12	0.02	0.31	0.31
	Sp <i>std3</i>	0.80	0.82	0.55	0.32	0.81
Gene level	MG (%) <i>std1</i>	62.7	69.7	95.3	27.9	41.8
	WG (%) <i>std3</i>	12.9	14.1	0.0	44.3	7.4

Table legends

Table Participating groups and associated annotation categories

Table Genefinding submissions

Table Evaluation of genefinding systems. Evaluation is divided into categories: Base level, Δ value, Gene level, Difference, statistical fit, ature, reported, sensitivity, **Sn**, Specificity, **Sp**, Missed (on **ME**), Wrong (on **WE**), Missed (e **MG**), Wrong (e **WG**), Sp (e **SG**), In (e **JG**). *std* and *std3* indicate gain which standard deviation is reported.

Table Evaluation of promoter prediction systems. Shows sensitivity for identified transcription start sites compared to false positive non-TSS regions and general gene region (all unlikely regions defined as here starting downstream from annotated transcription start (the energy regions spanning genome distance to previous annotated gene including the annotated TSS (Saknothe *std3* annotation)).

Table Evaluation of similarity search in Base gene level statistics show. The level described is sensitivity, **Sn**, specificity, **Sp**, the statistical gene level given Missed (e **MG**), Wrong (e **WG**).

References

- Agarwal, P., States, D. 1996. Comparative accuracy of methods for protein sequence similarity search. *Bioinformatics* **14**:40-47.
- Altschul, F., Gish, W., Miller, W., Myers, E., Lipman, D. 1990. A simple method for sequence search. *J. Mol. Biol.* **215**:403-410.
- Arkhipov, V. 1995. Promoter elements in *Drosophila melanogaster* revealed by sequence analysis. *Genetics* **139**:1359-1369.
- Ashburner, M. 2000. *submitted*.
- Ashburner, M. 1999. European *Drosophila* Genom Project (EDGP).
<http://edgp.ebi.ac.uk/>.
- Ashburner, M., Bor, R., Durbin, R., Guigo, R., Hubbard, T. 1999a. *GASP assessment meeting, EMBL Heidelberg*, .
- Ashburner, M., Misra, R., Root, E., Lewis, D., Blazek, D., David, C., Doyle, G., Gall, R., George, H., Harris, G., Hartze, D., Jarve, H., Hong, K., Houston, R., Hoskins, G., Johnson, M., Martin, M., Moshrefi, N., Palazzolo, M., Reese, S., Spradling, C., Sank, K., Wang, W., Whitelaw, K., Kimmahd, *et al.* 1999. An exploratory 20-Mb genomic region of *Drosophila melanogaster*. *Genetics* **153**:79-219.

Bateman A, Birney D, Durbin E, Eddy S, Howe E, Sonnhammer G. 2000 The Pfam Protein Family Database. *Nucleic Acids Res* **28**: 263-266.

Benson G. 1999 Tandem repeat finder programaly DNA sequence. *Nucleic Acids Res* **27**: 573-580.

Besemer H, Borodovsky M. 1999 Heuristic approach to deriving gene models. *Nucleic Acids Res* **27**: 11-3920.

Birney E. 1999 Wise2. <http://www.sanger.ac.uk/Software/Wise2/>.

Birney E, Durbin R. 1999 Dynamic flexible coding sequence comparison. *ISMB* **5**: 6-64.

Birney E, Durbin R. 2000 Using GeneWise to annotate a genome. *Genome Research* **10**: 1-10.

Burgin G, Karlin I. 1997 Prediction of complete gene structure from genomic DNA. *J Mol Biol* **268**: 8-94.

Burgin G, Karlin I. 1999 Finding genes in genomic DNA. *Curr Opin Struct Biol* **8**: 346-354.

Burset M, Guigó R. 1996 Evaluation of gene structure prediction programs. *Genomics* **34**: 353-367.

Cavaliere, P., Tuncel, G., Bonnard, B., Buche, J. 1999. The eukaryote promoter Database (EPD): recent developments. *Nucleic Acids Res* **27**: 307-309.

Cavaliere, P., Tuncel, G., Bonnard, B., Buche, J. 2000. The eukaryote promoter Database (EPD). *Nucleic Acids Res* **28**: 302-303.

Dunbrack, R.L., Gerloff, M., Bowler, C., Lichtarge, D., Eichen 1997. Meeting views on secondary structure prediction: Assessment of techniques for protein structure prediction (CASP2), San Jose, California, December 13-16, 1996. *Foldes* **2**: 27-42.

Eeckman, F.D., Durbin, R. 1995. A C-Base Database. *Methods Mol Biol* **48**: 583-605.

Fickel, W., Hatzigeorgiou, I. 1999. Eukaryote promoter recognition. *Genom Res* **7**: 861-878.

Fickel, W., Lung, J. 1999. Assessment of protein binding measures. *Nucleic Acids Res* **27**: 6441-6450.

Field, D. 1999. *unpublished.*

Flores, J., Hartzel, Z., Zhang, S., Rubin, W., Miller 1998. Computer program for alignment of DNA sequence with genomic DNA sequence. *Genom Res* **8**: 967-974.

Friesen, R., and Rubin, J. 1999.

*Proceedings of the Annual International
Conference on Computational Molecular Biology
(RECOMB), Lyon, France* ,

Gaasterland, T., Olesen, J., and MAGPIE. 1999. Automated genome interpretation.
Genet **126**:78.

ation. *Trends*

Green, P. 1995. *unpublished.*

Guigó, R., and Nudsen, D. 1999.

Smith, J. 1999. Prediction of gene structure.
M Biol
226:41-157.

M Biol

Harris, N., He, M., and Lewis, J. 1999. CloneCurator.

<http://www.fruitfly.org/displays/CloneCurator.html>.

Helmus, J. 1999. Neomorphic Genomes Software Development Toolkit

(NGSDK).

Neomorphic Inc, Berkeley. <http://www.neomorphic.com>.

Henikoff, S., Henikoff, J., and Petrokovski, J. 1999. New features of

the Block Database

servers. *Nucleic Acids Res* **27**:226-228.

Henikoff, S., and Henikoff, J. 1994. Protein family classification

searching

database. *Genomics* **19**:7-107.

Henikoff, S., and Henikoff, J. 1994b.

*27th Hawaii International Conference on System Sciences,
Hawaii, U.S.A.* ,

Henikoff, S., Henikoff, O. 2000. Genomic sequence annotation based on predicted protein-coding regions. *Genome Research* 10: 99-109.

Henikoff, S., Henikoff, O., Petrokovski, I. 1999. Blocks: a non-redundant database of protein domain architectures derived from multiple sequence alignments. *Bioinformatics* 15: 471-479.

Hubbard, T. 2000. Personal communication.

Jurkiewicz, K. 1998. Repeating DNA: mining and cleaning. *Curr Opin Struct Biol* 8: 333-337.

Krogh, A. 1997. A method for improving performance of HMMs in protein classification. *Ismb* 5: 179-186.

Kuhlmann, M., Haussler, M., Reese, M., Eckman, J. 1997. Integrating probabilistic structure models. *Proc Symp Biocomput* : 232-244.

Kurtz, S., Schleiermacher, S. 1998. REPEAT: a computational approach to repeats in complete genomes. *Bioinformatics* 14: 426-427.

Levit, M. 1997. Competitive assessment of fold recognition and alignment accuracy. *Proteins Suppl* 2: 92-104.

Marcotte, E.M., Pellegrini, M., Thompson, J., Yeates, T.L., Eisenberg, D. 1999. A combined algorithm for genome-wide prediction of protein function. *Nature* 402: 83-86.

Reese, *et al.* 11/28/2000 64

Mohr, 1999 EST_GENOME Programs for identifying DNA sequences unspliced genomic
DNA. *Comp Appl Biosci* **13**:477-478.

Moullier, Hubbar, Bryan, Fidalis, Pedersen, 1999 Critical assessment
method of protein structure prediction (CASP) round III. *Proteins Suppl*:6.

Moullier, Hubbar, Fidalis, Pedersen, 1999 Critical assessment
protein structure prediction (CASP): round III. *Proteins
Suppl*:6.

Ohler, Harbeck, Niemann, North, Rees, 1999 Interpolated
forkaryotic promoter recognition. *Bioinformatics* **15**:
362-369.

Ohler, Stommar, Harbeck, 2000 Stochastic Segmentation of
Promoter Regions. *Proc Symp Biocomput* **5**:377-388.

Parré, Blanc, Guigo, 2000 Gene ID in *Drosophila*. *Genome Research* **10**.

Pearson, 1995 Comparison of methods for searching protein sequence data
Sci **4**:1145-1160.

Pearson, Lipman, 1988 Improved fast algorithm for sequence comparison. *Proc
Natl Acad Sci* **85**:444-2448.

Reese, 2000 Genome annotation *Drosophila melanogaster* Ph.D. University of

Reese *et al.*
Hohenheim.
11/28/2000

Reese, M., Heckman, K., and Haussler, J. 1997. Improved splice site detection. *Genie. Comp Biol* **4**:1-323.

Reese, M., Harrig, G., Artz, H., Lewis, J. 1999. *The conference intelligent System Molecular Biology (SMB'99) Heidelberg, Germany* <http://www.fruitfly.org/GASP1>.

Reese, M., Kulp, D., Amman, A., Haussler, J. 2000. Gene finding in *Drosophila melanogaster*. *Genom Research* **10**.

Rubin, G. 2000. Full-length DNA project.

Rubin, G. 1999. Berkeley Drosophila Genom Project (BDGP). <http://www.fruitfly.org>.

Salamon, A., Solovyev, V. 2000. Gene finding in Drosophila genom DNA. *Genom Research* **10**.

Schen, M., Klingenhoff, W., Werner, J. 2000. *in preparation*.

Sipp, P., Ackner, D., Dominguez, S., Koppstein, J. 1999. *in preparation*. *Proteins Suppl* **2**:6-230.

Solovyev, V., Salamov, A., Lawrence, J. 1997. Identification of gene structures using discriminative functions and dynamic programming. *Ismb* **3**:67-375.

Sonnhammer E, Eddy R, Birney S, Bateman A, Durbin R 1998 family multiple
sequence alignments as HMM-profiles of
domains. *Nucleic Acids Res* 26:20-322.

Sonnhammer ES, Eddy R, Durbin R 1997 comprehensive database of protein
domain families based on alignments. *Proteins* 28:
405-420.

Stein D, Thierry-Mieg J 1998 Scriptable access to Caenorhabditis elegans
sequence in the ACEDB databases. *Genome Res* 8:
1308-1315.

Storm G 2000. *submitted.*

Uberbacher F, Burlingame 1991 Locating protein-coding regions in a
DNA
sequence by multiple sensor-neural network approach.
Protein Sci 8:1261-11265.

Zemla A, Venclova M, Moult J, Fidelis 1999 Processing signal
peptides by CASP3
protein structure predictions. *Proteins Suppl* 2-29.