# GenomeAnnotationAssessmentin *Drosophila melanogaster*

Martin G. Reese*, George Hartzell*, Nomi L. Harris*, Uwe Ohler* [†] and Suzanna E. Lewis*.

*Berkeley Drosophila Genome Project, Department of Molecular and Cel l Biology, UniversityofCalifornia,Berkeley,California94720-3200
[†]ChairforPatternRecognition,UniversityofErlangen-Nuremberg,Ma rtensstr.3,D-91058 Erlangen,Germany.


Sendproofsto:          MartinReese

                                    LawrenceBerkeleyNationalLaboratory

                                    1CyclotronRoad,MS64-121

                                    Berkeley,CA94720

                                    Phone:(510)486-4800

                                    Fax:(510)486-6798

                                    E-mail:mgreese@lbl.gov

# Abstract

Computationalmethodsforautomatedgenomeannotationarecriticaltoourcommunity'sa        bilityto
makefulluseofthelargevolumeofgenomicsequencebeinggeneratedandreleased.T        oexplore
theaccuracyoftheseautomatedfeaturepredictiontoolsinthegenomesofhigherorga        nismswe
evaluatedtheirperformanceonalarge,well-characterizedsequencecontigf        romthe *Adh*regionof
*Drosophilamelanogaster* .Thisexperiment,knownastheGenomeAnnotationAssessmentProject
(GASP),waslaunchedinMay1999.Twelvegroups,applyingstateofthearttools,contribute        d
predictionsforfeaturesincludinggenestructure,proteinhomologies,promotersites        ,andrepeat
elements.Weevaluatedthesepredictionsusingtwostandards,onebasedonpreviouslyunre        leased
highqualityfull-lengthcDNAsequencesandasecondbasedonthesetofannotationsgenerat        edas
partofanin-depthstudyoftheregionbyagroupofDrosophilaexperts(Ashburner        *etal.* ,1999b).
Whilethesestandardssetsonlyapproximatetheunknowndistributionoffeaturesinthisr        egion,we
believethatwhentakenincontexttheresultsofanevaluationbasedonthemaremeani        ngful.The
resultswerepresentedasatutorialattheconferenceonIntelligentSystem        sinMolecularBiology
(ISMB-99)inAugust1999(Reese     *etal.* ,1999).Over95percentofthecodingnucleotidesinthe
regionwerecorrectlyidentifiedbythemajorityofthegenefindersandthecorre        ctintrons/exon
structureswerepredictedformorethan40percentofthegenes.Homologybasedannotation
techniquesrecognizedandassociatedfunctionswithalmosthalfofthegenesinthere        gion,the
remainderwereonlyidentifiedbythe     *abinitio* techniques.Thisexperimentalsopresentsthefirst
assessmentofpromoterpredictiontechniquesforasignificantnumberofgenesinala        rge
contiguousregion.Wediscoveredthatthepromoterpredictors'highfalsepositiver        atesmaketheir
predictionsdifficulttouse.IntegratinggenefindingandcDNA/ESTalignmentsw        ithpromoter
predictionsdecreasesthenumberoffalsepositiveclassificati        onsbutdiscoverslessthanone-thirdof
thepromotersintheregion.Webelievethatbyestablishingstandardsforevaluati        nggenomic
annotationsandbyassessingtheperformanceofexistingautomatedgenomeannotationtools        ,this

experimentestablishesabaselinewhichcontributestothevalueofongoinglarge-s        caleannotation

projectsandshouldguidefurtherresearchingenomeinformatics.

# 1. Introduction:TheGenomeAnnotationAssessmentP roject(GASP)

Genomeannotationisarapidlyevolvingfieldingenomicsmadepossiblebythelarge-               scale
generationofgenomicsequencesanddrivenpredominantlybycomputationaltools.Thegoaloft            he
annotationprocessistoassignasmuchinformationaspossibletotherawsequenceofcompl             ete
genomeswithanemphasisonthelocationandstructureofthegenes.Thiscanbeaccomplishe            dby
*abinitio* genefinding,byidentifyinghomologiestoknowngenesfromotherorganisms,bythe
alignmentoffull-lengthorpartialmRNAsequencestothegenomicDNA,orthroughcombi            nations
ofsuchmethods.Relatedtechniquescanalsobeusedtoidentifyother           features,suchasthelocation
ofregulatoryelementsorrepetitivesequenceelements.Theultimategoalof           genomeannotation,
thefunctionalclassificationofalltheidentifiedgenes,currentl       ydependsondiscoveringhomologies
togeneswithknownfunctions.

Weareinterestedinanobjectiveassessmentofthestateof          theartinautomatedtoolsandtechniques
forannotatingcompletegenomes.TheGASPprojectwasorganizedtoformulateguidel          inesand
accuracystandardsforevaluatingcomputationaltoolsandtoencouragethedevelopment           ofnew
modelsandtheimprovementofexistingapproachesthroughacarefulassessmentandcompa           rison
ofthepredictionsmadebycurrentstate-of-the-artprograms.

TheGASPexperiment,thefirstofitskind,wassimilarinmanywaystotheCASP             (Critical
AssessmentofTechniquesforproteinstructureprediction)contestsforproteinstr         uctureprediction
(Dunbrack *etal.* ,1997;Levitt,1997;Moult    *etal.* ,1997;Moult   *etal.* ,1999;Sippl  *etal.* ,1999;
Zemla *etal.* ,1999),describedathttp://predictioncenter.llnl.gov.However,unliketheCASP
contest,GASPwaspromotedasacollaborationtoevaluatevarioustechniquesforgenom            e
annotation.

TheGASPexperimentconsistedofthefollowingstages:

- Trainingdataforthe *Adh*region,including2.9megabasesof *Drosophilamelanogaster* genomicsequence,wascollectedbytheorganizersandprovidedtotheparticipants.

- Asetofstandardswasdevelopedtoevaluatesubmissionswhiletheparticipating groups producedandsubmittedtheirannotationsfortheregion.

- Theparticipatinggroups'predictionswerecomparedtothestandards,ateamof independentassessorsevaluatedtheresultsofthecomparison,andtheresultswere presentedasatutorialatISMB-99.

Participantsweregiventhefinishedsequenceforthe *Adh*regionandsomerelatedtrainingdata,but theydidnothaveaccesstothefull-lengthcDNAsequencesthatweresequencedforthe paperby Ashburner *etal.* (1999b)thatdescribesthe *Adh*regionindepth.Theexperimentwaswidely announcedandopentoanyparticipants.Submitterswereallowedtouseanyavailable technologies andwereencouragedtodisclosetheirmethods.Sincewewerefortunatetoattracta largegroupof participantswhoprovidedawidevarietyofannotations,webelievethatourevaluati onaddresses thestateofartingenomeannotation.

TwelvegroupsparticipatedinGASP,submittingannotationsinoneormoreofsixcategor ies: *ab initio*genefinding,promoterrecognition,EST/cDNAalignment,proteinsimilarity,repe titive sequenceidentificationandgenefunction.Table1listseachparticipatinggroup,thenam esofthe programsorsystemsitused,andwhichofthesixclassesofannotationsitsubmitted.Addi tional papersinthisissuearewrittenbytheparticipantsthemselvesanddesc ribetheirmethodsandresults indetail.

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 in Ashburner *etal.* (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 dataset. At worst, there is a bias in the available data that makes our conclusions significantly misleading. We believe that the data is not unreasonable 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 dataset (see Birney and Durbin (this issue (2000)) for a further detailed discussion of evaluating these predictions).

In the next section we describe the target genomic sequence and the auxiliary data, including a critical discussion of our standard sets. Section 3 gives a short description of existing annotation methods that complement other papers in this issue, including a review article of existing gene finding methods by Stormo (2000) and papers describing the methods used by the individual participants. The Results section assesses the individual annotation methods and the Conclusion discusses what the experiment revealed about issues involved in annotating complete genomes. An article by Ashburner (2000) in this issue provides a biological perspective on the experiment.

## 2. Data: The benchmark sequence: The *Adh* region in *Drosophila melanogaster*

The selection of a genomic target region for assessing the accuracy of computational genome annotation methods was a difficult task for several reasons: The genomic region had to be large enough, the organism had to be well studied, and enough auxiliary data had to be available to have a

goodexperimentallyverified"correctanswer"butthedatashouldbeanonymoussothatabl       indtest

wouldbepossible.The   *Adh*regionofthe   *Drosophilamelanogaster* genomemetthesecriteria.

*Drosophilamelanogaster* isoneofthemostimportantmodelorganismsandalthoughthe       *Adh*

regionhadbeenextensivelystudied,thebestgeneannotationsandcDNAsfortheregionwer       enot

publisheduntilaftertheconclusionoftheGASPexperiment.The2.9megabase       *Adh* contigwas

largeenoughtobechallenging,containedgeneswithavarietyofsizesandstructure       s,andincluded

regionsofhighandlowgenedensity.Itwasnotacompletelyblindtest,however,sincese       veral

cDNAandgenomicsequencesforknowngenesintheregionwereavailablepriortothe

experiment.

## 2.1 GenomicDNAsequence

Thecontiguousgenomicsequenceofthe       *Adh* regioninthe   *Drosophilamelanogaster* genomespans

nearly3megabasesandhasbeensequencedfromaseriesofoverlappingP1andBACclonesa       sa

partoftheBerkeleyDrosophilaGenomeProject(Rubin&al.,1999)andtheEuropeanDrosophil       a

GenomeProject(Ashburner&al.,1999).Thissequenceisbelievedtobeofveryhighquality       with

anestimatederrorrateoflessthan1in10,000bases,basedonPHRAPqualityscores.Adeta       iled

analysisofthisregioncanbeaccessedthroughtheBDGPwebsite

(http://www.fruitfly.org/publications/Adh.html)aswellasinAshburner       *etal.* (1999b).

## 2.2 Curatedtrainingsequences

Weprovidedseveral   *Drosophilamelanogaster* specificdatasetstotheGASPparticipants.This

enabledparticipantstotunetheirtoolsforDrosophilaandfacilitatedacomparisonoft       hevarious

approachesthatwasunbiasedbyorganismspecificfactors.Thefollowingcuratedse       quencesets,

extractedfromFlybaseandEMBL,providedbytheEuropeanDrosophilaGenomeProjectat

Cambridge,andprovidedbytheBerkeleyDrosophilaGenomeProjectweremadeavail ableandcan

befoundathttp://www.fruitfly.org/GASP/data/data.html:

- Asetofcompletecodingsequences(starttostopcodon),excludingtransposableelements ,

  pseudogenes,non-codingRNAs,mitochondrialandviralsequences(2,122entries);

- Non-redundantsetofrepetitivesequences,notincludingtransposableele ments(96entries);

- Transposonsequences,containingonlythelongestsequenceofeachtransposonfamilyand

  excludingdefectivetransposableelements(44entries);

- GenomicDNAdatafrom275multi-and141single-exonnon-redundantgenestogether

  withtheirstartandstopcodonsandsplicesites,takenfromGenBankversion109;

- Asetof256unrelatedpromoterregions,takenfromEPD(CavinPerier *etal.* ,1999;Cavin

  Périer *etal.* ,2000)andacollectionmadebyI.Arkhipova(1995);

- AnuncuratedsetofcDNAandESTsequencesfromworkinprogressattheBerkeley

  DrosophilaGenomeProject.

Fiveoutofthetwelveparticipatinggroupsreportedmakinguseofthesedatasets .

## 2.3 Resourcesforassessingpredictions:The"correct"answer

Inacomparativestudythegoldstandardusedtoevaluatesolutionsisthemostimporta ntfactorin

determiningtheusefulnessofthestudy'sresults.Fortheresultstobe meaningful,thestandardmust

beappropriateandcorrectintheeyesofthestudy'saudience.Sinceourgoalwastoeval uatetools

thatpredictgenesandgenestructureincomplexeukaryoticorganismswedrewourst andardfroma

complexeukaryoticmodelorganism,choosingtoworkwitha2.9megabasesequencecontigfr om the *Adh*regionof *Drosophilamelanogaster* .Comparingpredictedannotationsinsucharegionis onlyconsequentialifthestandardisbelievedtobecorrect,ifthatcorrectnesshas beenestablished bytechniquesthatareindependentoftheapproachesbeingstudied,andifthepredictorshadno priorknowledgeofthestandard.Ideallyitwouldcontainthecorrectstructureofallt hegenesinthe regionwithoutanyextraneousannotations.Unfortunately,suchasetisimpossibletoobtains ince theunderlyingbiologyisincompletelyunderstood.Webuiltatwo-partapproximationtothe perfectdataset,takingadvantageofdatafromtheBDGPcDNAsequencingprojec t (http://www.fruitfly.org/EST)andaDrosophilacommunityefforttobuildasetofc urated annotationsforthisregion(Ashburner *etal.* ,1999b).Ourfirstcomponent,knownasthe *std1*data set,usedhighqualitysequencefromasetof80full-lengthcDNAclonesfromthe *Adh*regionto provideastandardwithannotationsthatareverylikelytobecorrectbutcertainl yarenot exhaustive.Thesecondcomponent,knownasthe *std3*dataset,wasbuiltfromtheannotations beingdevelopedforAshburner *etal.* (1999b)togiveastandardwithmorecompletecoverageofthe region,althoughwithlessconfidenceabouttheaccuracyandindependenceoftheannotations.We believethatthistwo-partapproximationallowsustodrawusefulconclusionsaboutthea bilityto accuratelypredictgenestructureincomplexeukaryoticorganismseventhoughthe absolutely perfectdatasetdoesnotexist.

Eukaryotictranscriptannotationshavecomplexstructuresbasedonthecomposition offundamental featuressuchastheTATAboxandothertranscriptionfactorbindingsites,thetranscr iptionstart site(TSS),thestartcodon,5-primeand3-primesplicesiteboundaries,thestopcodon,t hepoly-adenylationsignal,exonstartandendpositions,andcodingexonstartandendpositions.Ourgene predictionevaluationsfocusedonannotationsthatarespecifictothecodingregion,fromt hestart codonthroughthevariousintron-exonboundariestothestopcodon,andonpromoterannotations. Whileothertypesoffeaturesarealsobiologicallyinterestingwewereunabl etodevisereliable

methods for evaluating their predictions. Whenever possible we relied on unambiguous biological evidence for our evaluations; when that was not available we combined several types of evidence curated by domain experts.

Our goal for our first standard set, called *std1*, was to build a set of annotations that we believed were very likely to be correct in their fine details (e.g. exact locations for splice sites), even if we were unable to include every gene in the region. We based *std1* on alignments of 80 high quality, full-length cDNA sequences from this region with the high quality genomic sequence for the contig. The cDNA sequences are the product of a large cDNA sequencing project at the Berkeley Drosophila Genome Project and had not been submitted to GenBank at the time of the experiment. Working from five cDNA libraries, the longest clone for each unique transcript was selected and sequenced to a high quality level. Starting with these cDNA sequences, we generated alignments to the genomic sequence using sim4 (Florea *et al.*, 1998) and filtered them on several criteria. Of the eighty candidate cDNA sequences, three were paralogs of genes in the *Adh* region and nineteen appeared to be cloning artifacts (unspliced RNA or multiple inserts into the cloning vector), leaving us with alignments for fifty-eight cDNA clones. These alignments were further filtered based on splice site quality. We required that all of the proposed splice sites include a simple "GT"/"AG" core for the 5' and 3' splice sites respectively and that they scored highly (5' splice sites >= 0.35 threshold which gives a 98% true positive rate, and 3' splice sites >= 0.25 which gives a 92% true positive rate) using a neural network splice site predictor trained on *Drosophila melanogaster* data (Reese *et al.*, 1997). This process left us with forty-three sequences from the *Adh* region for which we had structures confirmed by alignments of high quality cDNA sequence data with high quality genomic data and by the fit of their splice sites to a Drosophila splice site model. Of these forty-three sequences, seven had a single coding exon and thirty-six had multiple coding exons. We added start codon and stop codon annotations to these structures from the corresponding records in the *std3* dataset.

After the experiment we recently discovered four inconsistent genes in the *std1* dataset. For two genes ( *DS07721.1, DS003192.4* ) the cDNA clones (CK02594, CK01083 respectively) are likely to be untranscribed genomic DNA that was inappropriately included in the cDNA library. Two other genes from *std3* (*DS00797.5* and *wb*) were incorrectly reported in *std1* as three partial all incomplete EST alignments (cDNA clones: CK01017, LD33192, and CK02229). In keeping with *std1*'s goal of highly reliable annotations, all four sequences have been removed from the *std1* data set that is currently available on the GASP website. The results reported here use the larger, less reliable, dataset as presented at the ISMB99 tutorial.

The complete set of the original 80 aligned high quality, full-length cDNA sequences was named *std2*. This set was never used in the evaluation process because it did not add any further compelling information or conclusions due to the unreliable alignments.

Our goal for the second, used standard set, called *std3*, was to build the most complete set of annotations possible while maintaining some confidence about their correctness. Ashburner *etal.* (1999b) compiled an exhaustive and carefully curated set of annotations for this region of the Drosophila genome based on information from a number of sources, included BLASTN, BLASTP (Altschul *etal.* ,1990), and PFAM alignments (Bateman *etal.* ,2000; Sonnhammer *etal.* ,1998; Sonnhammer *etal.* ,1997), high scoring GENSCAN (Burge & Karlin, 1997) and Genefinder (Green, 1995) predictions, ORFFinder results (Friese *etal.* ,1999), full length cDNA clone alignments (including those used in *std1*), and alignments with full length genes from GenBank. This set included 222 gene structures: 39 with a single coding exon, and 183 with multiple coding exons. Of these 222 gene structures, 182 are similar to a homologous protein in another organism or have a Drosophila EST hit. For these structures, the intron-exon boundaries were verified by partial cDNA/EST alignments using sim4 (Florea *etal.* ,1998), homologies were discovered using BLASTX, TBLASTX and PFAM alignments, and gene structure was verified using a version of

GENSCAN trained for finding human genes. Of the fifty-four remaining genes, fourteen had EST or homology evidence but were not predicted by GENSCAN or Genefinder, and forty were based entirely on strong GENSCAN and Genefinder predictions. All of this evidence was evaluated and edited by experienced Drosophila biologists, resulting in a protein coding gene dataset that exhaustively covers the region with a high degree of confidence and represents their view of what should or should not be considered an annotated gene. Their gene dataset excluded the seventeen found transposable elements (6 LINE-like elements (*G,F,Doc,* and *jockey*) and 11 retrotransposons with long terminal repeats (LTRs; *copia,roo,297,blood,mdg1-* like and *yoyo*), which almost all contain long ORFs. Some of these ORFs code for known and some others for so far unknown protein sequences.

Both of these datasets have shortcomings. As mentioned above, *std1* only includes a subset of the genes in the region. It also includes a pair of transcripts that represent alternatively spliced products of a single gene. While this is not incorrect, it confounds our scoring process. Because the cDNA alignments do not provide any evidence for the location of the start and stop codons, we based those annotations in *std1* on information from the *std3* set. Many of the gene structures in *std3* are based on GENSCAN and Genefinder predictions without other supporting evidence, so it is possible that the fine details are incorrect, that the entries are not entirely independent of the techniques used by the predictors in the experiment, and that the set overestimates the number of genes in the region.

See Birney and Durbin (this issue (2000)) and Henikoff and Henikoff (this issue (2000) for further discussion of the difficulties of evaluating these predictions especially in the protein homology annotation category, in which by training these programs will recognize protein like sequences such as the ORFs in transposable elements as genes. They and others (see other GASP publications in this issue) have raised the issues of annotation oversights, transposons, and pseudogenes. In cases where GASP submissions suggest a missed annotation this information has been passed onto

biologistsforfurtherresearch,includingscreeningcDNAlibraries.Webelie vethatitwouldhave beenbiasedtoretroactivelychangethescoringschemeusedattheGASPexperime ntbasedsolely onmissedannotationsdiscoveredbytheparticipant'ssubmissions.Seesection5foranexa mpleof anannotationthatmaybemissinginthestandarddatasets.Inthe *std3* datasetwebasedour standardforwhatisorisnotaDrosophilageneontheexpertannotationsprovidedin(Ashburner *et al.*,1999b).Itisclearthatbothtransposonsandpseudogenesaregenuinefeaturesofthegenome andthatgenefindingtechnologiesmightrecognizethem.Sincetheywerenotincludedas coding genesintheexpertannotations,wedecidedagainstincludingthemeinthestandardset.

Buildingasetfortheevaluationoftranscriptionstartsiteor,moregenerally ,forpromoter recognition,provedtobeevenmoredifficult.Forthegenesinthe *Adh*regionalmostno experimentallyconfirmedannotationforthetranscriptionstartsiteexists.A sthe5'UTRregionsin Drosophilacanextenduptoseveralkilobases,wecouldnotsimplyusetheregiondire ctlyupstream ofthestartcodon.Toobtainthebestpossibleapproximation,wetookthe5'endsofannotations fromAshburner *etal.* (1999b)wheretheupstreamregionreliedonexperimentalevidence(the5' endsoffull-lengthcDNAs)andforwhichthealignmentofthecDNAtothegenomicsequenc e includedagoodopenreadingframe.Theresultingsetcontained92genesoutofthe222 annotationsinthe *std3*set(Ashburner *etal.* ,1999b).Thisnumberislargerthanthenumberof cDNAsusedfortheconstructionofthe *std1*setdescribedabovebecauseweincludedcDNAsthat werealreadypubliclyavailable.The5'UTRofthese96geneshasanaverageleng thof1,860base pairs,aminimumlengthof0basepairs(whenthestartcodonwasannotatedatthebeginning,due tothelackofanyfurthercDNAalignmentinformation;thisisverylikelyt obeonlyapartial5' UTRandthereforeanannotationerror)andamaximumlengthof36,392basepairs.

## 2.4 Dataexchangeformat

Oneofthechallengesofageneannotationstudyisfindingacommonformatinwhichtoexpress thevariousgroups' predictions.Theformatmustbesimpleenoughthatallofthegroupinvol ved canadapttheirsoftwaretouseitandstillberichenoughtoexpressthevariousannotati ons.

WefoundthattheGeneralFeatureFormat(GFF)(formerlyknownastheGeneFeat ureFinding format)wasanexcellentfittoourneeds.TheGFFformatisanextensionofasimple< *name,start, end*>recordthatincludessomeadditionalinformationaboutthesequencebeingannotated:the sourceofthefeature;thetypeoffeature;thelocationofthefeatureinthesequence ;andascore, strand,andframeforthefeature.Ithasanoptionalninthfieldthatcanbeusedtogroupmultipl e predictionsintosingleannotations.MoreinformationcanbefoundattheGFFwebsite: http://www.sanger.ac.uk/Software/formats/GFF/.OurevaluationtoolsusedaGFFparserforthe PERLprogramminglanguagethatisalsoavailableattheGFFwebsite.

Wefoundthatitwasnecessarytospecifyastandardsetoffea turenameswithintheGFFformat,for instancedeclaringthatsubmittersshoulddescribecodingexonswiththefeaturename "CDS".We producedasmallsetofexamplefiles(accessiblefromtheGASPwebsite)tha twedistributedtothe submittersandwerepleasedwithhoweasilywewereabletoworkwiththeirres ults.

# 3. Methods

Genomeannotationisanongoingefforttoassignfunctionalfeaturestolocationsonthegenom ic DNAsequence.Traditionallymostoftheseannotationsrecordinformationaboutanorganis m's genes,includingproteincodingregions,RNAgenes,promotersandothergeneregulatoryel ements, aswellasgenefunction.Inadditiontothesegenefeatures,thefollowinggeneralge nomestructure

featuresarealsocommonlyannotated:repetitiveelementsandgeneralA,C,G ,Tcontentmeasures (e.g.,isochores).

## 3.1 Genomeannotationclasses

WhiletheGASPexperimentinvitedandencouragedanyclassofannotations,mostsubmiss ions wereforgene-relatedfeatures,emphasizing *abinitio* genepredictionsandpromoterpredictions.In addition,twogroupssubmittedfunctionalproteindomainannotationsandtwogroupssubmitted repeatelementannotations.Inthesectionsthatfollowwecategorizeanddiscusst hesubmitted predictions.

## 3.1.1 Genefinding

Proteincodingregionidentificationisamajorfocusofcomputationalbiology.Asepara tearticlein thisissue(Stormo,2000)discussesandcomparescurrentmethods,whileanearlypaperby Fickett andTung(1992)andamorerecentreviewofgeneidentificationsystemsbyBurgeandK arlin (1998)giveexcellentoverviewsofthefield.Table2liststhesixgroupsthatpredic tedprotein-codingregionswiththecorrespondingprogramnames.Italsocategorizesthesubmiss ionsbasedon thetypesofinformationusedtobuildthemodelforpredictions.Whileallgroupsusedstatis tical informationfortheirmodels-predominantlycodingbias,codingpreference,andconsensus sequencesforstartcodon,splicesitesandstopcodons-onlytwogroupsusedproteinsimilari ty informationorpromoterinformationtopredictgenestructure.Morethanhalfofthegroups incorporatedsequenceinformationfromcDNAsequences.Ingeneral,state-of-the-a rtgene predictionsystemsusecomplexmodelsthatintegratemultiplegenefeatures intoaunifiedmodel.

## 3.1.2 Promoterprediction

Thecomplicatednatureofthetranscriptioninitiationprocessmakescomputational        promoter
recognitionahardproblem.Wedefinepromoterpredictionastheidentificationoftransc        ription
startsites(TSS)ofproteincodinggenesthataretranscribedbyeukaryoticR        NApolymeraseII.A
detaileddescriptionofthestructureofpromoterregionsandexistingpromoter        predictionsystemsis
beyondthescopeofthispaper.FickettandHatzigeorgiou(1997)providea        nexcellentreviewofthe
field.

Wecanbroadlyidentifythreedifferentapproachestopromoterprediction,withatlea        stoneGASP
submissionineachcategory.Thefirstclassconsistsof"searchbysi        gnal"programs,whichidentify
singlebindingsitesofproteinsinvolvedintranscriptioninitiation,orcombinationsofs        itesto
improvethespecificity.TheprogramCoreInspectorbyWerner'sgroup(Scherf        *etal.*,2000)
belongstothiscategoryandsearchesforco-occurrencesoftwocommonbindingsitesw        ithinthe
corepromoter(thecorepromoterusuallydenotestheregionwherethedirectcontactbe        tweenthe
transcriptionmachinery,theholoenzymeofthetranscriptioncomplex,andtheDNAtake        splace).
Thesecondclassisoftentermed"searchbycontent",asprogramswithinthisgroupdonotr        elyon
specificsignalsbuttakethemoregeneralapproachofidentifyingthepromoterre        gionasawhole,
frequentlybasedonstatisticalmeasures.Sometimesthepromoterissplitint        oseveralregionsto
obtainmoreaccuratestatistics.TheMCPromoterprogram(Ohler        *etal.*,1999)isamemberofthis
secondgroup.Incomparisonwiththesignal-basedgroup,thecontent-basedsystemsusually        are
moresensitivebutlessspecific.Thethirdclasscanbedescribedas"promoterpr        edictionthrough
genefinding".Simplyusingthestartofagenepredictionasaputativetranscript        ionstartsitecanbe
verysuccessfulifthe5'UTRregionisnottoolarge.Thisapproachcanbeimprovedbyinc        luding
homologytoESTsequencesand/orapromotermoduleinthestatisticalsystemsusedfor        gene

prediction.TheTSSpredictionssubmittedbytheparticipantsoftheMAGPIEandthe        Geniegroups

belongtothislastclass.


Thenotoriousdifficultyoftheproblemitselfisexacerbatedbythelimitedamountofe        xisting

reliablyannotatedtrainingmaterial.TheexperimentalmappingofaTSSisala        boriousprocessand

isthereforenotroutinelycarriedout,evenifthegeneitselfisstudiedextensi        vely.So,bothtraining

themodelsandevaluatingtheresultsisadifficulttask,andtheconclusionswedraw        fromthe

resultsmustbeconsideredwithmuchcaution.


### 3.1.3 Repeatfinders

Detectingrepeatedelementsplaysaveryimportantroleinmodelingthe3-dim        ensionalstructureof

aDNAmolecule,specificallythepackingoftheDNAinthecellnucleus.Itisbel        ievedthatthe

packingoftheDNAaroundthenucleosomeiscorrelatedwiththeglobalsequencestructur        e

producedpredominantlybyrepetitiveelements.Repeatsalsoplayamajorroleinev        olution(fora

reviewsee(Jurka,1998)).Twogroups,GaryBenson(TandemRepeatsFinderversion2.02

(TRF)(Benson,1999))andtheMAGPIEteamusingtwoprograms(Calypso(Field,)andR        EPuter

(Kurtz&Schleiermacher,1999)submittedrepetitivesequenceannotations.TRF(Be        nson1999)

locatesapproximatetandemrepeats(i.e.,twoormorecontiguous,approximatecopiesofapa        ttern

ofnucleotides)wherethepatternsizeisunspecifiedbutfallswithintherangefr        om1to500bases.

TheCalypsoprogram(Field,)isanevolutionarygenomicsprogram.Itsprimaryf        unctionistofind

repetitiveregionsinDNAandproteinsequencesthathavehigherthanaveragemuta        tionrates.The

REPuterprogram(Kurtz&Schleiermacher,1999)determinesrepeatsofafixedpr        e-selectedlength

incompletegenomes.

### 3.1.4 Proteinhomologyannotation

Homologiestogenesequencesfromotherorganismscanoftenbeusedtoidentifyprotein-          coding
regionsinanonymousgenomicsequence.Inadditiontothelocation,itisoftenpossibletoinfert          he
functionofthepredictedgenebasedonthefunctionofthehomologousgeneintheotherorganism
orofaknownstructuralandfunctionalproteinelementinthegene.Whilethetoolsinthegene
predictioncategoryandtheEST/cDNAalignmentcategoryareusuallyintendedt          odeterminethe
exactstructureofagene,theproteinhomologybasedtoolsareusuallyoptimizedtofindcons          erved
partsofthesequencewithoutworryingabouttheexactgenestructure.Traditionallyt          hisareaof
genomeannotationshasbeendominatedbythesuiteoflocalalignmentsearchtoolsofBLA          ST
(Altschul *etal.* ,1990)andmoreglobalsearchtoolssuchasFASTA(Pearson&Lipman,1988).
Recentreviewsinthisareainclude(Agarwal&States,1998;Marcotte          *etal.* ,1999;Pearson,1995).

IntheGASPexperimenttwogroupsspecializinginfunctionalproteindomainormotif
identificationingenomicDNAsubmittedannotations.TheHenikoffgroupfoundhitstothe
BLOCKS+database(http://blocks.fhcrc.org),adatabaseconsistingofconserve          dproteinmotifs
(Henikoff *etal.* ,1999a;Henikoff&Henikoff,1994a;Henikoff&Henikoff,1994b).Thesecond
groupinthiscategorysubmittedresultsfromtheGeneWiseprogram(Birney,1999).T          hisprogram
searchesgenomicDNAagainstacomprehensiveHMM-basedlibrary(PFAM,(          Bateman *etal.* ,
2000;Sonnhammer *etal.* ,1998;Sonnhammer *etal.* ,1997))ofproteindomains.Bothprogramslook
forconservedregionsbysearchingtranslatedDNAagainstarepresentationofm          ultiplealigned
sequences.WhileinBLOCKS+themultipleproteinalignmentsconsistofsetsofunga          ppedregions
theGeneWiseprogramsearchesagainstagappedalignment.Bothmethodswillturnupdi          stantly
relatedsequences.

## 3.1.5 EST/cDNA alignment

Computational predictions of gene location and structure go hand in hand with EST/cDNA sequencing and alignment techniques for building transcript annotations in genomic sequence. Either can be used as a discovery tool, with the other held in reserve for verification. A researcher can verify the existence and structure of predicted genes by sequencing the corresponding mRNA molecules and aligning their sequences to the original genomic sequence. Alternatively, one can start with an EST or cDNA sequence and build an alignment to the genomic sequence that has been guided and/or verified by tools from the gene prediction arsenal; for example using likely splice site locations, and checking for long open reading frames and potential frameshifts.

There are many tools for aligning sequences. While they have generally been specialized for aligning sequences that are evolutionarily related, some are designed for niche applications such as recognizing overlaps among sequencing runs. Aligning EST/cDNA sequences to the original genomic sequence also presents a unique set of tradeoffs and issues. In some cases (inter-species EST/genomic alignments) these tools must model evolutionary changes in the sequence. Sometimes (e.g. for low quality EST sequences) they need to model errors in the sequence generated by the sequencing process. For multi-exon genes, they need to model the intron regions as cost-free gaps tied to a model for recognizing splice sites. Several tools have been developed for this task: Mott (1997) and Birney and Durbin (1997) describe dynamic programming approaches that include models of splice sites and intron gaps. Florea *et al.* (1998) describe *sim4*, a heuristic tool that performs as well as the dynamic programming approaches and is efficient enough to support searching of large databases of genomic sequence.

Using cDNA clones and their sequences to build transcript annotations requires a variety of operations. The tools discussed above align the cDNA sequences to the genomic sequence, but steps must be taken to filter out clones that are merely paralogs of genes in the sequence and to

recognize and handle various laboratory artifacts. If the clones represent short ESTs, then a likely annotation can be built by assembling a consistent model from their individual alignments. Longer ESTs or cDNAs might generate several similar alignments, and an automated tool must be able to select the most biologically meaningful variant. While there are some gene prediction tools that can use information about homologies to known genes or ESTs, and most large scale sequencing centers have some automated sanity checking for their database search results, there are not any tools that automate the production of transcript annotations from cDNA sequences.

### 3.1.6 Gene function

Gene function predictions are the most difficult annotations to produce and to evaluate. Current technologies use similarity to proteins (or protein domains) with known function to predict functional domains in genomic sequence. While some tools use simple sequence alignments, more powerful tools have developed significantly more sensitive models.

It quickly became apparent that a consistent and correct assessment of function predictions as part of the GASP experiment was not possible due to the incomplete understanding of the protein products encoded by the 222 genes in the *Adh* region.

### 3.2 Evaluating gene predictions

An ideal gene prediction tool would produce annotations that were exactly correct and entirely complete. The fact that no existing tool has these characteristics reflects our incomplete understanding of the underlying biology as well as the difficulty to build adequate gene models in a computer. While no tool is perfect, each tool has particular strengths and weaknesses and any performance evaluation should be in the context of an intended use. For example, researchers who are interested in identifying gene rich regions of a genome for sequencing would be happy with a

toolthatsuccessfullyrecognizesagene'sapproximatelocation,evenifitincor rectlydescribed
splicesiteboundaries.Ontheotherhand,someonetryingtopredictproteinstructuresism ore
interestedingettingagene'sstructureexactlyrightthaninatool'sabilit ytopredicteverygenein
thegenome.

Whenassessingtheaccuracyofpredictions,eachpredictionfallsi ntooneoffourcategories.Atrue
positive(TP)predictionisonethatcorrectlypredictsthepresenceofafeatur e.Afalsepositive
(FP)predictionincorrectlypredictsthepresenceofafeature.Atruenegati ve(TN)predictionis
correctinnotpredictingthepresenceofafeaturewhenitisn'tthere.Afalseneg ative(FN)
predictionfailstopredicttheexistenceofafeaturethatactuallyexists.T hesensitivity(Sn)ofa
toolisdefinedasTP/(TP+FN),andcanbethoughtofasameasureofhowsuccessfulthetool isat
findingthingsthatarereallythere.Thespecificity(Sp)ofatoolisdefinedas TP/(TP+FP),andcan
bethoughtofasameasureofhowcarefulatoolisaboutnotpredictingthingsthataren'tr eally
there.BursetandGuigó(1996)alsouseacorrelationcoefficientandanaveragecorr elation
coefficient.Wechosenottousethesemeasuresbecausetheydependonpredictors'true negative
informationandwerecognizethatourevaluationsetswereconstructedinsuchawayt hatthetrue
negativeinformationisnottrustworthy.Thesesensitivityandspecificitym etricsareusedfor
evaluatingthesubmissionsinthegenefinding,promoterrecognitionandgeneidentifica tionusing
proteinhomologycategories.Inthegenefindingcategorytheyareusedforallthree levels:base
level,exonlevelandgenelevel.Intheproteinhomologycategorytheyareusedforbase leveland
genelevelonly.

Inoneofthefirstreviewsofgenepredictionaccuracy,Fickett andTung(1992)developedamethod
thatmeasuredpredictors'abilitytocorrectlyrecognizecodingregionsingenom icsequence.They
usedtheirmethodtocomparepublishedtechniquesandconcludedthatin-framehexamercounts
werethemostaccuratemeasureofaregion'scodingpotential.Bur setandGuigó(1996)recognized

thatthereareawidevarietyofusesforgenepredictionsanddevelopedmeasures           --includingbase

level,exonlevel,andgenelevelspecificityandsensitivity--thatdescr          ibeapredictor'ssuitabilityfor

aparticulartask.

## 3.2.1 Baselevel

Thebaselevelscoremeasureswhetherapredictorisabletocorrectlylabel          abaseinthegenomic

sequenceasbeingpartofsomegene.Itrewardspredictorsthatgetthebroadsweepsof          agene

correct,eveniftheydon'tgetthedetailssuchasthespliceboundariesentir          elycorrect.It

penalizespredictorsthatmissasignificantportionofthecodingsequence,evenift          heygetthe

detailscorrectforthegenestheydopredict.Weusedthesensitivityandspecifi          citymeasures

definedaboveasthemeasuresofsuccessinthiscategory.

## 3.2.2 Exonlevel

Exonlevelscoresmeasurewhetherapredictorisabletoidentify          exonsandcorrectlyrecognizetheir

boundaries.Beingoffbyasinglebaseateitherendoftheexonmakesthepredictionincorre          ct.

Sinceweonlyconsideredcodingexonsinourassessment,thefirstexonisbracketedbythes          tart

codonanda5'splicesite,thelastexonisbracketedbya3'splicesiteandthestopcodon,andthe

interiorexonsarebracketedbyapairofsplicesites.Asmeasuresofsuccessi          nthiscategory,we

usedtwostatisticsinadditiontosensitivityandspecificity.The          *missedexon* (ME)scoreisa

measureofhowfrequentlyapredictorcompletelyfailedtoidentifyanexon(nopredicti          onoverlap

atall),whilethe    *wrongexon* (WE)scoreisameasureofhowfrequentlyapredictoridentifiesan

exonthathasnooverlapwithanyexoninthestandardsets.TheMEscoreisthepercentageof

exonsinthestandardsetforwhichtherewerenooverlappingexonsinthepredictedset.Simi          larly,

theWEscoreisthepercentageofexonsinthepredictedsetforwhichtherewerenoover lapping exonsinthestandardset.

### 3.2.3 Genelevel

Genelevelsensitivityandspecificitymeasurewhetherapredictorisabl etocorrectlyidentifyand assembleallofagene'sexons.Forapredictiontobecountedasatruepositive,alloft hecoding exonsmustbeidentified,everyintron-exonboundarymustbeexactlycorrect,andallofthee xons mustbeincludedinthepropergene.Thisisaverystrictmeasurethataddressesat ool'sabilityto perfectlyidentifyagene.Inadditiontothesensitivityandspecificitymeas uresbasedonabsolute accuracy,weusedthe *missedgenes* (MG)scoreasameasureofhowfrequentlyapredictor completelymissedagene(astandardgeneisconsideredmissedifnoneofitsexonsar eoverlapped byapredictedcodinggene)andthe *wronggenes* (WG)scoreasameasureofhowfrequentlya predictorincorrectlyidentifiedagene(apredictionisconsideredwrongifnoneofit sexonsare overlappedbyagenefromthestandardset).

### 3.2.4 SplitandJoinedgenes

Theexonlevelscoresdiscussedabovemeasurehowwellapredictorrecognizesexonsa ndgets theirboundariesexactlycorrect.Thegenelevelscoresmeasurehowwellapredi ctorcanrecognize exonsandassemblethemintocompletegenes.Neitherofthesescoresdirectlymea suresa predictor'stendencytoincorrectlyassembleasetofpredictedexonsintomoreorfe wergenesthan itshould.Wedevelopedtwonewmeasures, *splitgenes* (SG)and *joinedgenes* (JG) ,whichdescribe howfrequentlyapredictorincorrectlysplitsagene'sexonsintomultiplegenesa ndhowfrequently apredictorincorrectlyassemblesmultiplegenes'exonsintoasinglegene.Bec ausethecoverageof the *std1*datasetissoincomplete,wehaveonlyincludedsplitgenesandjoi nedgenescoresfromthe

comparisonwith *std3*.Agenefromthestandardsetisconsidered *split*ifitoverlapsmorethanone predictedgene.Similarly,apredictedgeneisconsidered *joined*ifitoverlapsmorethanonegenein thestandardset.TheSGmeasureisdefinedasthesumofthenumberofpredictedgenest hat overlapeachstandardgenedividedbythenumberofstandardgenesthatweresplit.Simi larly,the JGmeasureisthesumofthenumberofstandardgenesthatoverlapeachpredictedgenedi videdby thenumberofpredictedgenesthatwerejoined.Ascoreof1isperfectandmeansthatall ofthe genesfromonesetoverlapexactlyonegenefromtheotherset.


### 3.2.5 Applicationofthesemeasurestocorrectanswerdatasets *std1/std3*

Webuiltthe *std1*datasetinsuchawaythatwebelieveitiscorrectinthedetailsofthegenest hatit describes,thoughweknowthatitonlyincludesasmallportionofthegenesintheregion.The *std3* dataset,ontheotherhand,isascompleteaswaspossible,butdoesnothaverigorousindependent evidenceforallofitsannotations.Forthe *std1*dataset,webelievethattheTPcount(itwas predictedanditexistsinthestandard)andFNcount(itwasnotpredictedbutitdoesexisti nthe standard)arereliablebecauseoftheconfidencethatwehaveinthecorrectnessof thepredictionsin theset.Ontheotherhand,wedonotbelievethattheTNcount(itwasnotpredictedanditisnotin thestandardset)andFPcount(itwaspredictedbutisnotinthestandardset)arereli ablebecause theybothassumethatthestandardcorrectlydescribestheabsenceofafeatureandw eknowthat therearegenesmissingfrom *std1*.Itfollowsthatwebelievethatsensitivityismeaningfulfor *std1* becauseitonlydependsonTPandFNbutthatwearelessconfidentaboutthespecificitysc ore, sinceitdependsonTPandFP.Asimilarlogicappliestothe *std3*dataset,whereourconfidencein theset'scompletenessbutnotitsfinedetailssuggeststhattheTPandFPscore sareusablebutthat theTNandFNscoresarenot.Thismeansthatfor *std3*,webelievethatthespecificitymeasurecan beusedtodescribeapredictor'sperformancebutthatsensitivityislikelytobe misleading.

## 3.3 Evaluationofpromoterpredictions

WeadoptedthemeasuresproposedbyFickettandHatzigeorgiou(1997).Theyevaluatedt he successofpromoterpredictionsbygivingthepercentageofcorrectlyidentifie dtranscriptionstart sitesversusthefalsepositiverate.ATSSisregardedasidentifiedif aprogrammakesoneormore predictionswithinacertain"likely"regionaroundtheannotatedsite.Thefalsepos itiverateis definedasthenumberofpredictionswithinthe"unlikely"regionsoutsidethe"likely "regions dividedbythetotalnumberofbasescontainedintheunlikelyset.AsourannotationoftheTSSi s onlypreliminaryandnotexperimentallyconfirmed,wechosearatherlargeregionof 500bases upstreamand50basesdownstreamoftheannotatedTSSasthe"likely"region.Theupstr eam regionisalwaystakenasthe"likely"region,evenifitoverlapswithanei ghboringgeneannotation onthesamestrand.The"unlikely"regionforeachgenethenconsistsoftherestoftheg ene annotation,frombase51downstreamoftheTSStotheendofthefinalexon.

## 3.4 Visualizationoftheannotations

Generating"good"annotationsgenerallyrequiresintegratingmultiplesources ofinformation,such astheresultsofvarioussequenceanalysistoolsplussupportingbiologicalinformat ion. Visualizationtoolsthatdisplaysequenceannotationsinabrowsablegraphicalfram eworkmakethis processmuchmoreefficient.Inthisexperimentwefoundthatvisualizationtoolsare essentialin ordertoevaluatethegenomeannotationsubmissions.Whenannotationsaredisplayedvisua lly, overalltrendsbecomeapparent,forexamplegene-richvs.gene-poorregions;genest hatwere predictedbymostparticipantsvs.thosethatwerepredictedbyfew.Additionally,as wediscuss below,avisualizationtoolthatiscapableofdisplayingannotationsatmultiplelev elsofdetail providesawaytoexamineindividualpredictionsindetail.

Building genome annotation visualization tools is a daunting task. Many such tools have been developed, starting with ACeDB (Eeckman & Durbin, 1995; Stein & Thierry-Mieg, 1998). We were fortunate in that the Berkeley Drosophila Genome Project has built a flexible suite of genome visualization tools (Helt & al., 1999) that could be extended to display the GASP submissions. We adapted the BDGP's annotated clone display and editing tool, CloneCurator (Harris *et al.*, 1999) which is based on a genomic visualization toolkit (Helt & al., 1999), to read the annotation submissions in GFF format and display each team's predictions in a unique color and location.

CloneCurator (see Figure 1) displays features on a sequence as colored rectangles. Features on the forward strand appear above the axis, while those on the reverse strand appear below the axis. The display can be zoomed and scrolled to view areas of interest in more detail. A configuration file identifies the feature types that are to be displayed, and assigns colors and offsets to each one. For example, the *std1* and *std3* exons appear in yellow and orange close to the central axis.

## 4. Genome annotation results

The results of an experiment such as GASP are only meaningful if enough groups participate. We were fortunate to have twelve diverse groups involved and we were very grateful for the speed with which they were able to submit their predictions. We believe that these twelve groups provide a fair representation of the state of the art in annotation system technology. We collected submissions by electronic mail and evaluated them using the *std1* and *std3* datasets as described above. Before releasing our results at the Intelligent Systems in Molecular Biology conference in August 1999 in Heidelberg, Germany, we assembled a team of independent assessors (Ashburner *et al.*, 1999a) to review our techniques and conclusions. As discussed in the introduction, the accuracy of the various measures discussed below depends heavily on how well our standard sets capture the true

setoffeaturesintheregion.Thesevaluesshouldonlybeconsideredinthecontextofthest andard datasets.

Adetaileddescriptionoftheresultsandtheevaluationtechniquesweusedcanbeacces sedthrough theGASPhomepageat http://www.fruitfly.org/GASP/.

## 4.1.1 Genefinding

Table3summarizestheperformanceofthegenefindingtoolsusingthemeasuresdefi nedabove. Threegroupssubmittedmultiplesubmissions.Thefirstgroup,Fgenes1-3,submittedthre e predictionsatvaryingstringency(fordetailssee(Salamov&Solovyev,2000) ).FortheGeneID program,twosubmittedversionsarepresented,version1(GeneIDv1)beingtheorigina l submissionandversion2(GeneIDv2)beinganewersubmissionfromacorrectedversionof the originalprogram(fordetailssee(Parra *etal.*,2000)).Thethirdgroupwithmultiplesubmissions usedthreeversionsoftheGenieprogram:thefirstapurestatisticalapproach( Genie),thesecond includingESTalignmentinformation(GenieEST)andthethirdusingprot einhomologyinformation (GenieESTHOM)(fordetailssee(Reese *etal.*,2000)).ForallothergroupsfromTable2onlyone submissionwasevaluated.Thefollowingsectionsdiscussthebaselevel,exonle vel,andgenelevel performanceofthesesubmissions.

## 4.1.1.1 Baselevelresults

Severalgenepredictiontoolshadasensitivityofgreaterthan0.95atthebaselev el.Thissuggests thatcurrenttechnologyisabletocorrectlyidentifyover95%ofthe *Drosophilamelanogaster* proteome.Afewtoolsdemonstratedaspecificityofgreaterthan0.90atthebaselev el,only infrequentlylabelinganon-codingbaseascoding.Generallythetoolshaveahigherse nsitivity

thanspecificity.Twoprograms,Fgenes2andGeneID,weredesignedtobeconservati         veabouttheir

predictionsanddonotfollowthistrend.


## 4.1.1.2 Exonlevelresults

Therewasagreatdealofvariabilityintheexonlevelscores.Severaltools         hadsensitivityscores

around0.75,correctlyidentifyingbothexonboundariesabout75%ofthetime.Theirspecificitie         s

weregenerallymuchlower(thehighestwas0.68),probablyareflectionofthestric         tdefinitionof

exonlevelscoresbothsplicesiteshadtobepredictedcorrectly-a         ndpossibleinaccuraciesinthe *std3*

dataset.ThelowMEscores(severalbelow0.05)combinedwiththef         airlyhighsensitivitiessuggest

thatseveraltoolsweresuccessfulatidentifyingexonsbuthadtroublefindingthecor         rectexon

boundaries.ProgramsthatincorporateESTalignmentinformation,suchasGenieESTa         nd

HMMGene,hadsensitivityscoresthatwereupto10%betterthantheothertools.ThehighW         E

scoressuggesteitherthatthetoolsareover-predictingorthattherearege         nesthataremissingeven

from *std3*.


## 4.1.1.3 Genelevelresults

Allofthepredictorshadconsiderabledifficultycorrectlyassemblingcomplet         egenes.Thebest

toolswereabletoachievesensitivitiesbetween0.33and0.44,meaningthattheyareinc         orrecta

littleoverhalfofthetime.Thisvalueseemstobeverysimilarin         *Drosophila* andhuman

sequences,basedonarecentanalysisofthe         *BRCA2*regioninhuman(Hubbard,2000).Evenonthe

morecomplete *std3*dataset,theprogramstendedtoincorrectlypredictmanygenes.Theverylow

MGscore(aslowas4.6%)isreassuringsinceitsuggeststhatseveraltool         sareabletorecognizea

gene,eveniftheyhavedifficultyfiguringouttheexactdetailsofitsstructur         e.ComparingtheWG

andMGmeasuressuggeststhatexistingtoolstendtopredictgenesthatdonot         existmoreoftenthan

theymissgenesthatdoexist. Sinceitisalmostcertainthattherearereal genesthataremissing frombothstandardsets, thisconclusionmustbeviewedwithsomeskepticism. Whilethere were severaltoolswithgoodSGorJGscores, noneofthemperformedwellinbothcategories.

## 4.1.2 Promoterprediction

Table4showstheperformanceofthepromoterpredictionsystems, groupedbyapproach: search-by-signal, search-by-region, andgenepredictionprograms.

GenefindingprogramsthatincludeapredictionoftheTSSobtainedthebestresults. Thenumber offalsepredictionsmadebytheregion-basedprogramsisveryhigh(givingthem alowspecificity), andsincethesignalspecificprogramsonlyidentifyonepromotertheirsensitivityisverylow. The highspecificityofthegenefindersisobviouslyduetothecontextinformation: allpromoter predictionswithingenepredictionsareruledoutinadvance, andthelocationofthepossiblestart codonprovidesthesystemwithagoodinitialguessofwheretolookforapromoter. TheMAGPIE systemalsousesESTalignmentstoobtaininformationon5'UTRs, whichmirrorsthewaythe *std* setswereconstructed: roughlyonethirdoftheputativeTSSassignmentsrelyoncDNAsthatwere publiclyavailableinGenBank. Acloserlookattheresultsrevealsthattheregionbasedprograms haveasensitivitythatiscomparabletothegenefindersandthesignalbasedprogramhadonlya singlefalsepositive, showingthatbothtypesoftoolscanbeusedfordifferentapplications.

Ourdataset, andtheevaluationbasedonit, reliesontheassumptionthatthe5'endsofthefull-lengthcDNAsarereasonablyclosetothetranscriptionstartsite. Thismakesitveryhardtodraw strongconclusionsfromthepresentedresults. Eventhemostsensitivesystemscouldidentifyonly roughlyonethirdofthestartsites. Thiscouldofcoursebecausedbythefactthattheexisting annotationisonlyanapproximationandsomeofthetruetranscriptionstartsitesmaybelocated furtherupstream. Italsohintsatthediversityofpromoterregionsthatmirrors thepossibilitiesfor

generegulation,andattheexistingbiastowardshousekeepinggenesinthecurrentdata        setsused

forthetrainingofthemodels.


### 4.1.3 Geneidentificationusingproteinhomology

Genefindingevaluationstatistics,suchasthosedescribedinsection4.1.1,canbeusedto

summarizetheabilityofaprogramtoidentifycompleteandaccurategenestruc        turesingenomic

DNA.InTable5wehaveappliedthesameevaluationstatisticstothehomologybase        dsearch

programsGeneWiseandBLOCKS+.Becausetheseprogramsarenotoptimizedtodeal        withexact

exonboundaryassignments,Table5onlyshowstheperformanceforthebaselevelandthemiss        ed

andwronggenes.


Theverylowsensitivitiesatthebaselevelarenotsurprising,becausethepr        ogramsidentifyonly

conservedproteinmotifsorparticulardomainsandmakenoefforttopredictcompleteg        enes.

Specificity,whichshouldbehighgiventhatonlyconservedproteinmotifsarescored,w        aslower

thanexpected.Detailedstudiesofthesepredictions(see(Birney&Durbin,2000;Henik        off&

Henikoff,2000)inthisissue)showthatmostofthefalsepositivepredictionswerehit        sto

transposableelementsortogenesthataremissinginthestandardsets.Bothprogr        amsusea

databaseofproteindomainsorconservedproteinmotifs.Bothdatabasesarelargeandar        ebelieved

tocontainatleast50%oftheexistingproteindomains.ThehighnumberofMG,62.7%for

BLOCKSand69.7%forGeneWise,meansthattheseprogramswillmissasignifica        ntnumberof

DrosophilageneswhenusedtosearchgenomicDNAdirectly.TheWGscoresof12.9%BLOCK        S

and14.1%forGeneWisearelowerthanthegenefindingprogramsdiscussedintheprevious

section.

## 4.1.4 GeneidentificationusingEST/cDNAalignments

ItisbelievedthatsomecDNAinformationexistsforapproximatelyhalfofthege nesinthe *Drosophilamelanogaster* genome.ThiscDNAdatabase(availableastheESTdatasetattheGASP website)wasusedasabasisforthecDNA/ESTalignmentcategory.Thesensi tivityof31%for MAGPIEESTandGrailSimilarity(Table5)implythatthecodingportionoftheav ailableESTdata currentlycoversonethirdofthegenome'scodingsequence.Thelowspecificityisv erysurprising andsuggeststhattheEST/cDNAalignmentproblemisnotatrivialone.Theonly programthattried toaligncompletecDNAstogenomicDNA,MAGPIEcDNA,couldfindcompletecDNAsf oronly 2.4%ofthegenes.ESTalignmentsalsoresultedinhighnumbersofmissedgenes,sugges tingthat theESTlibrariesarebiasedtowardshighlyexpressedgenes.ThehighW Gscoressuggestthatsome genesaremissingevenfrom *std3*.

## 4.1.5 Selectedgeneannotations

Thesummarystatisticsdiscussedaboveonlyprovideaglobalviewofthepredicti ngprograms characteristics.Amuchbetterunderstandingofhowthevarious014approa chesbehavecanbeobtained bylookingatindividualgeneannotations.Suchadetailedexaminationcanal sohelpidentifyissues thatarenotaddressedbycurrentsystems.

Inthefollowingparagraphs,wewilldiscussafewinterestingexamples.Figure 1showsthecolor codesoftheparticipatinggroupsthatareusedthroughoutthissection.Geneslocatedonthe topof eachmaparetranscribedfromdistaltoproximal(withrespecttothetelomereof chromosomearm 2L);thoseonthebottomaretranscribedfromproximaltodistal. *Std1*and *std3* aretheexpert annotationsdescribedinAshburner *etal.* (1999b).Justbelowtheaxis,youcanseetheannotations forthetworepeatfindingprograms.Thesehavenosequenceorientationandaretherefore only shownononeside.Fartherawayfromtheaxis,after *std1*and *std3*,wegroupedallofthe *abinitio*

gene-findingprogramstogether.Nexttothegenefindersarethehomology-basedannotat        ions.On

thebottomandthetopofthefigureweshowthethreepromoterannotations,butforclaritywedi        d

notincludetheseannotationsinthesubsequentfigures.(Onthefrontpageandinthelegendof

Figure1,youcanseethefullsetofannotationsofallprograms,whicharealsoaccess        iblefromthe

GASPwebsite.)

Ourfirstexampleisa"busy"regionwithtwelvecompletegenesandonepartialg        eneinastretchof

onlyfortykilobases(Figure2A).Thisregionislocatedatthe3'endofthe        *Adh*regionfrombase

2,735,000tobase2,775,000.Genesexistonbothstrandsanditisstrikingthatinthisregionthe

genestendtoalternatebetweentheforwardandthereversestrands.Weselect        edthisregionforits

genedensityandbecauseithascharacteristicsthataretypicalofthecomple        te *Adh*region.Figure

2Avividlydemonstratesthatallofthegene-findingprograms'predictionsarehi        ghlycorrelated

withtheannotatedgenesfrom        *std1/std3*.Inthepastgenefindershadoftenmistakenlypredicteda

geneonthenon-codingstrandoppositeofarealgene,leadingtofalse        positivepredictionsknownas

"shadowexons".Figure2Amakesitclearthatgenefindershaveovercomethisproble        m,since

therearealmostnoshadowexonpredictionsforanyofthegenesin        *std3*.Anothercharacteristic,

capturedinthehighbaselevelsensitivityandthelowmissinggenesstatisti        cs,isthateverygenein

the *std3*setwaspredictedbyatleastafewgroupsandthatmostofthesepredictionsagre        ewith

eachother.Exceptforthesecondandthirdgenes(        *DS02740.5,I(2)35Fb* )ontheforwardstrand

(2,740,000-2,745,000),whichseemtobesingleexongenes,allofthegenesinthisregionare

multi-exongeneswithbetweentwoandeightexons.Theexonsizevarieswidely.Therear        egenes

thatconsistofonlytwolargeexons,somethatconsistofamixoflargeandsmallexons,ands        ome

thataremadeupexclusivelyofmanysmallexons.Thedistributionseemstobealmostr        andom.

Exceptforthelongfinalintroninthelastgeneonthereversestrand(        *cact*),theregionconsists

exclusivelyofshortintrons.

Predictions on the reverse strand indicate a possible gene from base 2,741,000 to base 2,745,000. Most of the gene finders agree on this prediction but neither *std1* nor *std3* describes a gene at this location. This could be a real gene that was missed by the expert annotation pathway described in Ashburner *et al* (1999b). Neither BLOCKS+ nor GeneWise found any homologies in this region, but we can see from the table in the previous section that many real genes do not have any homology annotations. Interestingly, this is the only area in the region where two gene finders predicted a possible gene that likely consists of shadow exons.

The fifth gene on the forward strand (*DS02740.10*, bases 2,752,500-2,755,000) shows that long genes with multiple exons are much harder to predict than single exon gene or genes with only a few exons. In this region splitting and joining genes does not seem to be a problem. Repeats occur sparsely and mostly in non-coding regions, predominantly in introns.

In contrast to the "busy" region in Figure 2A, Figure 2B highlights a region of almost equal size in which only one gene (*DS01759.1*) is present in both *std1* and *std3*. There are very few false positive predictions by any group, but there is one case where the "false" predictions by different programs are located at very similar positions (on the reverse strand near base 620,000). This suggests a real gene that is missing from both standard sets.

Figures 3A-3D depict selected genes that illustrate some interesting challenges in gene finding. Figure 3A show the *Adh* and the *Adhr* genes that occur as gene duplicates. The encoded proteins have a sequence identity of 33%. The positions of the two introns interrupting the coding regions are conserved and give additional evidence to tandem duplication. Both genes are under the control of the same regulatory promoter, the *Adhr* gene does not have a transcription start site of its own and its transcript is always found as part of an *Adh-Adhr* dicistronic mRNA. Gene duplications occur very frequently in the Drosophila genome-estimates show that at least 20% of all genes occur in gene family duplications. In an additional twist, *Adh* and *Adhr* are located within an intron

ofanothergene, *outspread*( *osp*),thatisfoundontheoppositestrand(fordetailsseeFigure3B). *Adh*iscorrectlypredictedbymostoftheprograms,althoughoneerroneouslypredictsanaddit ional firstexon.Mostoftheprogramsalsopredictthestructureof *Adhr*correctly;oneprogrammissesthe initialexonandshortensthesecondexon.Both *Adh*and *Adhr*showhitstotheproteinmotifsin BLOCKS+aswellasalignmentstoaPfamproteindomainfamilythroughGeneWi se.Bothgenes hittwodifferentPfamfamiliesandtheorderofthesetwodomainsisconservedinthe gene structure.

Figure3Bhighlightsthe *outspread*( *osp*)generegion.Thisisanexampleofagenewith exceptionallylong(>20kilobasepairs)introns,makingithardforanygenefindertopre dictthe entirestructurecorrectly.Inaddition,thereareanumberofsmallergenes(inc ludingthe *Adh*and *Adhr*genesdiscussedabove, *DS09219.1* (r.)and *DS07721.1*(f.))withintheintronsof *outspread*. Nocurrentgenefinderincludesoverlappinggenestructuresinitsmodel;asaconseque nce,noneof theGASPgenefinderswereabletopredictthe *outspread*structurewithoutdisruption.Thisis clearlyashortcomingoftheprogramssincegenescontainingothergenesareofte nobservedin *Drosophila*(Ashburner *etal.* reportseventeencasesforthe *Adh*region).However,itshouldbe notedthatmostofthegenefinderspredictthe3'endof *outspread* correctlyandthereforegetmost ofthecodingregionright.Theregionthatincludesthe5'endof *outspread*showsalotofgene predictionactivitybutthereisn'tanyconsistencyamongthepredictions.Oneprogra m (FGenesCCG3)doescorrectlypredictthe *DS09219.1* gene.

Figure3Cshowstheentiregenestructureofthe *Ca-alpha1D* gene.Thisgeneisthemostcomplex geneinthe *Adh*region,withmorethanthirtyexons.Thisisaverygoodexampleforstudyinggene splitting.Severalpredictorsbreakthegeneupintoseveralgenesbutsomegr oupsmakesurprisingly closepredictions.Thisshowsthecomplexstructurethatgenescanexhibitandthatextent towhich thiscomplexityhasbeencapturedinthestate-of-the-artpredictionmodels.Itisi nterestingtonote

thatmostofthelargerexonsarepredictedwhiletheshorterexonsaremissed.        Suchalargecomplex geneisagoodcandidateforalternativesplicing,whichcanultimately        bedetectedonlybyextensive cDNAsequencing.

Figure3Dshowsthetripleduplicationofthe        *idgf*gene(  *idgf1*, *idgf2*, *andidgf3* )ontheforward strand.Twoprogramsmistakenlyjointhefirsttwogenesintoasinglegene;all        theotherscorrectly predictallthreegenes.

# 5. Discussion

ThegoaloftheGASPexperimentwastoreviewandassessthesta        teoftheartingenomeannotation tools.Webelievethatthenoncompetitiveframeworkandthecommunity'senthusiast        ic participationhelpedusachievethatgoal.Byprovidingalloftheparticipantswitha        n unprecedentedsetof  *D.melanogaster* trainingdataandusingunreleasedinformationaboutthe regionasourgoldstandard,wewereabletoestablishthelevelplayingfieldthatm        adeitpossibleto comparetheperformanceofthevarioustechniques.Thelargesizeofthe        *Adh*contigandthe diversityofitsgenestructuresprovideduswithanopportunitytocomparethecapabili        tiesofthe annotationtoolsinasettingthatmodelsthegenomewideannotationscurrentlybeingatt        empted. However,thelackofacompletelycorrectstandardsetmeansthatourresultss        houldonlybe consideredinthecontextofthe        *std1*and *std3*datasets.

## 5.1 Assessingtheresults

Themostdifficultpartoftheassessmentwasdevelopingabenchma        rkforthepredictedannotations. Bydividingthepredictionsintodifferentclassesanddevelopingclass-specifi        cmetricsthatwere basedonthebestavailablestandards,wefeelthatwewereabletomakeameaning        fulevaluationof

thesubmissions.Whilemostoftheinformationthatwasusedtoevaluatethesubmissions was unreleased,somecDNAsequencesfromtheregionwereinthepublicdatabases.Asseque ncing projectsmoveforward,itwillbecomeincreasinglydifficultforfutureexperi mentstofindsimilarly unexploredregions.ThismakesitverydifferentfromtheCASPproteinstructurepr ediction contests,whichcanusethe3-dimensionalstructureofanoveltargetproteinthatis unknowntothe predictors.

Asdiscussedintheintroduction,thelackofanabsolutelycorrectst andardagainstwhichtoevaluate thevariouspredictionsisatroublingissue.Whilewebelievethatthestandardset ssufficiently representthetruenatureoftheregionandthatconclusionsbasedonthemareinteresting ,itmustbe rememberedthatthevariousresultscanonlybeevaluatedinthecontextoftheseinc ompletedata sets.ThisalsomakesGASPmoredifficultandlessclearcutthanCASP,where the3-dimensional proteinstructureisexperimentalsolvedatleasttosomedegreeofresolution.

Itshouldalsobenotedthatthegenefindingtoolswiththehighestspecificityhaveag reatdealin commonwithGENSCAN,thegenepredictiontoolusedinthedevelopmentofthe *std3*dataset. Thissuggeststhat *std3*'soriginsmighthaveledtoabiasfavoringGENSCAN-likepredictors. Because *std1*wasexclusivelycreatedusingfull-lengthcDNAalignments,thissetmight bebiased towardshighlyexpressedgenes,becausethecDNAlibrarieswerenotnormalized.

## 5.2 Progressingenomewideannotation

Therapidreleaseofcompletedgenomes,includingtheimminentreleaseofthe *Drosophila melanogaster*andhumangenomes,hasdrivensignificantdevelopmentsingenomeannotationand genefindingtools.Problemsthathaveplaguedgenefindingprograms,suchaspredictings hadow exons,restrictingpredictionstoasinglestrand,recognizingrepeats,andaccura telyidentifying

splicesiteshavebeenovercomebythecurrentstateoftheart.Inthissection,we          discusssomeof theremainingissuesingenomeannotationthattheGASPexperimenthighlighted.

Successfulgenepredictionprogramsusecomplexmodelsthatintegrateinformati          onfromstatistical featuresthataredrivenbythe3-dimensionalprotein-DNA/RNAinteractions          .Theymakeintegrated predictionsonbothstrandsandhavebeentunedtopredictallthegenesingenerichregionsand avoidover-predictinggenesingenepoorregions(Figure2Aand2B).Whilemostoftheprogr          ams identifyalmostalltheexistinggenes(asevidencedbythesensitivityandmi          ssinggenestatistics) thereissignificantvariationintheirabilitytoaccuratelypredictpre          cisegenestructures(seethe specificitystatistics,particularlyattheexonlevel).Ifanyglobalper          formanceconclusioncanbe drawnitisthattheprobabilisticgenefinders(mostlyHMMbased)seemtobemore          reliable.The integrationofEST/cDNAsequenceinformationintothe       *abinitio* genefinders(seeHMMGene, GenieEST,andGRAIL(Figures2B-2F))significantlyimprovesgenepredicti          ons,particularlythe recognitionofintron-exonboundaries.Somegroupssubmittedmultipleannotationsofthe       *Adh* regionusingprogramsthatweretunedfordifferenttasks.ThesuiteofFgenesprogr          amsshowsvery nicelytheresultsofsucha3-partsubmission.ThefirstFgenessubmission(FGene          s1)isaversion adjustedtoweightsensitivityandspecificityequally.Thesecondsubmission(F       Genes2)isvery conservativeandonlyannotateshigh-scoringgenes.Thisresultsinahighspecifi          citybutalow sensitivity.Thethirdsubmission(FGenes3)triestomaximizesensitivity          andavoidmissingany genes,atthecostofalossinspecificity.Thesedifferentlytuned          variantsmaybeusefulfordifferent typesoftasks.

Acomparison(datanotshown)toagenefindingsystemthatwastrainedonhumandatashowed thatitdidnotperformaswellastheprogramsthatweretrainedon          *Drosophila*data.

Noneofthegenepredictorsscreenedfortransposableelements,whichhaveaprot          ein-likestructure. AsdescribedinAshburner     *etal.* (1999b),the *Adh*regionhasseventeentransposableelement

sequences.Eliminatingtransposonsfromthepredictionsoraddingthemtothestandardset swould havereducedthefalsepositivecounts,raisingthespecificityandloweringtheW EandWGscores. Whilethisaccountsforaportionofthehighfalsepositivescoreswebelievethat theremayalsobe additionalgenesinthisregionnotannotatedin *std3*.Futurebiologicalexperiments(Rubin,2000)to identifyandsequencethepredictedgenesthatwerenotincludedin *std3*shouldimprovethe completenessandaccuracyofthefinalannotations.

Therewerefewersubmissionsofhomology-basedannotationsthanthoseby *abinitio* genefinders andtheirresultsweresignificantlyaffectedbytheirfalsepositivera tes.Asignificantportionof thosefalsepositiveswerematchestotransposableelements,someappeartobem atchestopseudo-genes,andothersarelikelytoberealbutasyetun-annotatedgenes.Thehomology-base d approachesseemtobethemostpromisingtechniqueforinferringfunctionsfornewlypredi cted genes.

EvenusingEST/cDNAalignmentstopredictgenestructuresisnotas simpleasexpected.Paralogs, lowsequencequalityofmRNAs,andthedifficultyofcloninginfrequentlyexpressedmRN Asmake thismethodofgenefindingmorecomplexthanbelievedanditisdifficulttoguarantee completenesswiththismethod.NormalizedcDNAlibrariesandothermoresophistic ated technologiestopurifygeneswithlowexpressionlevels,alongwithimprovedalignme ntand annotationtechnologies,shouldimprovepredictionsbasedonEST/cDNAalignments.

## 5.3 LessonsfortheFuture

Inordertofullyassessthesubmittedannotations,the"correctanswer"mustbeimpr oved.Only extensivefull-lengthcDNAsequencingcanaccomplishthis.Apossibleapproachwouldbe to designprimersfrompredictedexons/genesinthegenomicsequenceandthenusehybridizat ion technologiestofishoutthecorrespondingcDNAfromcDNAlibraries.Forpromoterpredi ctions,

anotherwaytoimprovethe"correctanswer"istomakegenome-to-genomealig        nmentswiththe

DNAofrelatedspecies(e.g.,    *C.Briggsae* versus  *C.elegans* ; *D.melanogaster* versus  *D.virilis* ).

Moredetailedguidelines,includinghowtohandleambiguousfeaturessuchaspseudogenesand

transposons,willmaketheresultsoffutureexperimentsevenmoreuseful.

Asuccessfulsystemtoidentifyallgenesinagenomeshouldconsistofacombinationof        *abinitio*

genefinding,EST/cDNAalignments,proteinhomologymethods,promoterrecognitionandrepe        at

finding.Allofthevarioustechnologieshaveadvantagesanddisadvantagesandanautoma        ted

methodforintegratingtheirpredictionsseemsideal.

Beyondtheidentificationofgenestructureisthedeterminationofgenefunctions.Most        ofthe

existingprototypesofsuchsystemsarebasedonsequencehomologies.Whilethisis        agoodstarting

pointitisdefinitelynotsufficient.Thestateoftheartforpredictingfunctioninpr        oteinsequences

usestheprotein'sthree-dimensionalstructure,butthedifficultyofaccurately        predictingthree-

dimensionalstructurefromprimarysequencesmakesapplyingthesetechniquesoncom        plete

genomesproblematic.Thenewfieldofstructuralgenomicswillhopefullygivemor        eanswersin

theseareas.

Anotherapproachtofunctionclassificationistheanalysisofgenee        xpressiondata.Improvementsin

transcriptionstartsiteannotations,alongwithcorrelationinexpressionprofiles        ,shouldbevery

helpfulinidentifyingregulatoryregions.

## 6. Conclusions

TheGASPexperimentsucceededinprovidinganobjectiveassessmentofcurrentappr        oachesto

geneprediction.Themainconclusionsfromthisexperimentarethatcurrentmethodsofge        ne

predictionsaretremendouslyimprovedandthattheyareveryusefulforgenomescal        eannotations

butthathighqualityannotationsalsodependonasolidunderstandingoftheorganisminquestion (*e.g.,*recognizingandhandlingtransposons).

ExperimentslikeGASPareessentialforthecontinuedprogressofautomatedannota tionmethods. Theyprovidebenchmarkswithwhichnewtechnologiescanbeevaluatedandselected.

ThepredictionscollectedinGASPshowedthatformostofthegenesoverlappingpredict ionsfrom differentprogramsexisted.Whetherornotacombinationofoverlappingpredictionswoulddo betterthanthebestperformingindividualprogramwasnotexplicitlytestedinthis experiment.For suchatestadditionalexperimentssuchascDNAlibraryscreeningandsubsequentfull -length cDNAsequencinginthisselected *Adh* testbedregionwouldbenecessary.Theseexperimentsare currentlyunderwayanditwouldbeinterestingtoperformasecondGASPexperimentwhenm ore cDNAshavebeensequenced.

Webelievethatexistingautomatedannotationmethodsarescalableandthattheulti matetestwill occurwhenthecompletesequenceofthe *Drosophilamelanogaster* genomebecomesavailable. Thisexperimentwillsetstandardsfortheaccuracyofgenome-wideannotationandim provethe credibilityoftheannotationsdoneinotherregionsofthegenome.

# 7. URLs

## 7.1 Genefinding

**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

## 7.2 Promoterprediction

**MCPromoter:**                http://www5.informatik.uni-

erlangen.de/HTML/English/Research/Promoter

**CoreInspector:**             http://www.gsf.de/biodv

## 7.3 Proteinhomology

**BLOCKS+:**                   http://blocks.fhcrc.org

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

**GeneWise:**                  http://www.sanger.ac.uk/Software/Wise2/

## 7.4 Repeatfinders

**TRF:**                       http://c3.biomath.mssm.edu/trf.test.html
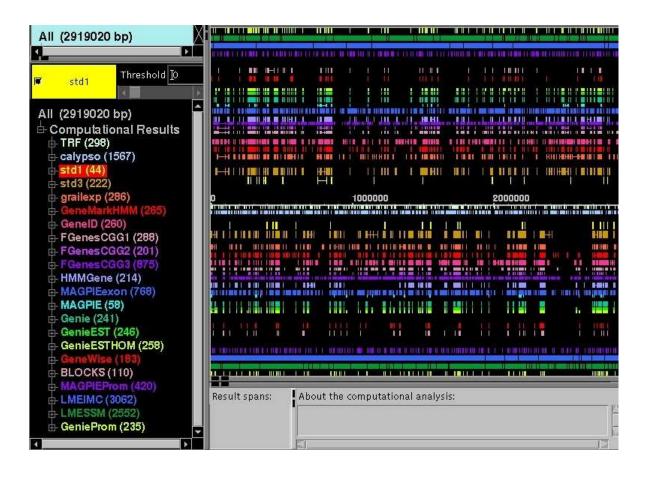
# 8. Acknowledgments

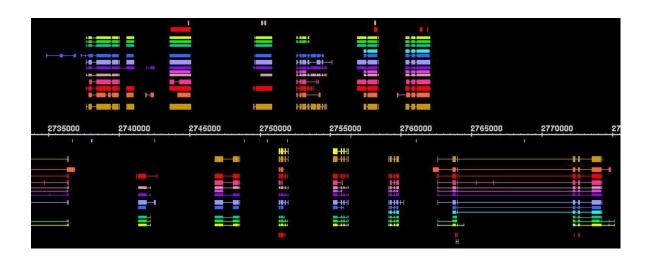# Figures

*Figure1(GASP)*

*Figure2A(busyregion)*

*Figure2B(desert)*

*Figure3A(Adh-Adhr)*

*Figure3B(outspread)*

*Figure3C(Ca-alpha1D)*

*Figure3D(idgf)*

# Figurelegends

## *Figure1(GASP)*

ThisfigureisascreenshotfromtheCloneCuratorprogram(Harris *etal.* ,1999).Itfeaturesthe genomeannotationsofall12groupsforthe2.9megabase *Adh*region.Themainpanelshowsthe computationalannotationsontheforward(aboveaxis)andreversesequencestrands(below axis). Geneslocatedonthetophalfofeachmaparetranscribedfromdistaltoproximal(withr espectto thetelomereofchromosomeare2L);thoseonthebottomaretranscribedfromproximaltodi stal. Rightbelowtheaxisarethetworepeatfindingresultsdisplayed,followedbyrefer encesetsfrom Ashburner *etal.* ( *std1*and *std3*),followedbythetwelvesubmissionsofgenefindingprograms, followedbythetwoproteinhomologyprogramsandeventually,farthestawayfromtheaxi s,the fourpromoterrecognitionprograms.Theleftpanelgivesthecolor-codedlegendforthepr ogram andthenumberofpredictionsmadebytheprograms.

| Program Identifier | Color | Referenceinthis GenomeResearch Issue | Reference |
|---|---|---|---|
| TRF | seafoam | | Benson(1999) |
| Calypso | lightblue | | Field(1999) |
| std1 | yellow | | Unpublishedconservativealignmentof cDNAs |
| std3 | orange | | Ashburner *etal.* (1999b) |
| Grailexp | redorange | | UberbacherandMural(1991) |
| GeneMarkHMM | red | | BesemerandBorodovsky(1999) |
| GeneID | hotpink | | Guigó(1992) |
| FGenesCGG1 | pink | | Solovyev *etal.* (1995) |
| FGenesCGG2 | magenta | | Solovyev *etal.* (1995) |
| FGenesCGG3 | purple | | Solovyev *etal.* (1995) |
| HMMGene | cornflower | | Krogh(1997) |
| MAGPIEexon | blue | | Gaasterland(1996) |
| MAGPIE | turquoise | | Gaasterland(1996) |
| Genie | seagreen | | Reese *etal.* (1997) |
| GenieEST | green | | Kulp(1997) |
| GenieESTHOM | chartreuse | | Kulp(1997) |
| GeneWise | red | | unpublished |
| BLOCKS | pink | | Henikoff *etal.* (1999b) |
| MAGPIEProm | purple | | unpublished |
| LMEIMC | blue | | Ohler *etal.* (1999) |
| LMESSM | darkgreen | | Ohler *etal.* (2000) |
| GenieProm | chartreuse | | Reese(2000) |

*Figure2A(busyregion)*

AnnotationsforthefollowingknowngenesdescribedinAshburner *etal.* (1999b)areshownforthe regionfrom2,735,000-2,775,000(fromthelefttotherightofthemap):

*crp*(partial,rev.), *DS02740.4*(f), *DS02740.5*(f), *I(2)35Fb*(f), *heix*(r), *DS02740.8*(f), *DS02740.9* (r), *DS02740.10*(f), *anon-35Fa*(r), *Sed5*(f), *cni*(r), *fzy*(f), *cact*(r).

*Figure2B(desert)*

Annotationsforthefollowingknowngenedescribedin Ashburner *etal.* areshownfortheregion from600,000-635,000(fromthelefttotherightofthemap):

*DS01759.1*(r).

*Figure3A(Adh-Adhr)*

Annotationsforthefollowingknowngenesdescribedin Ashburner *etal.* areshownfortheregion from1,109,500-1,112,500(forwardstrandonly)(fromthelefttotherightofthemap):

*Adh,Adhr* .

*Figure3B(outspread)*

Annotationsforthefollowingknowngenesdescribedin Ashburner *etal.* areshownfortheregion from1,090,000-1,180,000(fromthelefttotherightofthemap):

*outspread*or *osp*(r), *Adh*(f), *Adhr*(f), *DS09219.1*(r), *DS07721.1*(f).

*Figure3C(Ca-alpha1D)*

Annotationsforthefollowingknowngenedescribedin Ashburner *etal.* areshownfortheregion

from2,617,500-2,640,000(forwardstrandonly)(fromthelefttotherightofthemap):

*Ca-alpha1D.*


*Figure3D(idgf)*

Annotationsforthefollowingknowngenesdescribedin Ashburner *etal.* areshownfortheregion

from2,894,000-2,904,000(forwardstrandonly)(fromthelefttotherightofthemap):

*idgf1,idgf2,idgf3* .

# Tables

*Table1:ParticipatingGroupsandassociatedannotationcategories*

| | Programname | Gene finding | Promoter recognition | EST/cDNA Alignment | Protein Similarity | Repeat | Gene function |
|---|---|---|---|---|---|---|---|
| Mural *etal.* Oakridge,US | **GRAIL** | X | | X | | | X |
| Parra *etal.* Barcelona,ES | **GeneID** | X | | | | | |
| Krogh Copenhagen, DK | **HMMGene** | X | | | | | |
| Henikoff *etal.* Seattle,US | **BLOCKS** | | | | X | | X |
| Solovyev *etal.* Sanger,UK | **FGenes** | X | | | | | |
| Gaasterland *et al.* Rockefeller, US | **MAGPIE** | X | X | X | | X | X |
| Benson *etal.* MountSinai, US | **TRF** | | | | | X | |
| Werner *etal.* Munich, GER | **CoreInspector** | | X | | | | |
| Ohler *etal.* Nuremberg, GER | **MCPromoter** | | X | | | | |
| Birney Sanger,UK | **GeneWise** | | | | X | | X |
| Reese *etal.* Berkeley/SantaCruz,US | **Genie** | X | X | | | | |

*Table2:Genefindingsubmissions*

| | Program name | Statistics | Promoter | EST/cDNA Alignment | Protein similarity |
|---|---|---|---|---|---|
| **Mural *etal.* Oakridge,US** | **GRAIL** | **X** | | **X** | |
| **Guigó *etal.* Barcelona,ES** | **GeneID** | **X** | | | |
| **Krogh Copenhagen, DK** | **HMMGene** | **X** | | **X** | **X** |
| **Solovyev *etal.* Sanger,UK** | **FGenes** | **X** | | | |
| **Gaasterland *et al.* Rockefeller, US** | **MAGPIE** | **X** | **X** | **X** | |
| **Reese *etal.* Berkeley/Sant aCruz,US** | **Genie** | **X** | **X** | **X** | **X** |

*Table3*

| | | Fgenes 1 | Fgenes 2 | Fgenes 3 | Gene ID v1 | Gene ID v2 | Gen ie | Gen ie EST | Gen ie EST HOM | HMM Gen e | MAG PIE exo n | GRA IL |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Base level | Sn *std1* | 0.89 | 0.49 | 0.93 | 0.48 | 0.86 | 0.96 | 0.97 | 0.97 | 0.97 | 0.96 | 0.81 |
| | Sp *std3* | 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 *std1* | 0.65 | 0.44 | 0.75 | 0.27 | 0.58 | 0.70 | 0.77 | 0.79 | 0.68 | 0.63 | 0.42 |
| | Sp *std3* | 0.49 | 0.68 | 0.24 | 0.29 | 0.34 | 0.57 | 0.55 | 0.52 | 0.53 | 0.41 | 0.41 |
| | ME(%) *std1* | 10.5 | 45.5 | 5.6 | 54.4 | 21.1 | 8.1 | 4.8 | 3.2 | 4.8 | 12.1 | 24.3 |
| | WE(%) *std3* | 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 *std1* | 0.30 | 0.09 | 0.37 | 0.02 | 0.26 | 0.40 | 0.44 | 0.44 | 0.35 | 0.33 | 0.14 |
| | Sp *std3* | 0.27 | 0.18 | 0.10 | 0.05 | 0.10 | 0.29 | 0.28 | 0.26 | 0.30 | 0.21 | 0.12 |
| | MG(%) *std1* | 9.3 | 34.8 | 9.3 | 44.1 | 13.9 | 4.6 | 4.6 | 4.6 | 6.9 | 4.6 | 16.2 |
| | WG(%) *std3* | 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 |

*Table4*

| SystemName | Sensitivity | Rateoffalsepositive predictionsinregion(a) (853,180bases) | Rateofpredictionsin region(b)(2,570,232 bases) |
|---|---|---|---|
| CoreInspector | **1(1%)** | **1/853,180** | **1/514,046** |
| MCPromoterV1.1 | **26(28.2%)** | **1/2,633** | **1/2,537** |
| MCPromoterV2.0 | **31(33.6%)** | **1/2,437** | **1/2,323** |
| GeniePROM | **25(27.1%)** | **1/14,710** | **1/28,879** |
| GenieESTPROM | **30(32.6%)** | **1/16,729** | **1/29,542** |
| MAGPIE | **33(35.8%)** | **1/14,968** | **1/16,370** |

*Table5*

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

## Tablelegends

**Table1:** ParticipatingGroupsandassociatedannotationcategories

**Table2:** Genefindingsubmissions

**Table3:** Evaluationofgenefindingsystems.Theevaluationisdividedinthreecategories :Base level,exonlevelandGenelevel.Thedifferentstatisticalfeaturesreport edareSensitivity( **Sn**), Specificity( **Sp**),MissedExon( **ME**),WrongExon( **WE**),MissedGene( **MG**),WrongGene( **WG**), SplitGene( **SG**)andJoinedGene( **JG**)." *std1*"and" *std3*"indicateagainstwhichstandardsetthe statisticsarereported.

**Table4:** Evaluationofpromoterpredictionsystems.Weshowthesensitivityforidentifie d transcriptionstartsitesincomparisontothefalsepositiveratefornon-TSS regionsandgeneral generegions:(a)the"unlikely"regiondefinedastherestofthegenestarting 50basesdownstream fromitsannotatedtranscriptionstartsite;(b)thegeneralgeneregion,spanning fromhalfthe distancetothepreviousandnextannotatedgenesincludingtheannotatedTSS(takenfromthe *std3* annotation).

**Table5:** Evaluationofsimilaritysearching.Baseandgenelevelstati sticsareshown.Thebaselevel isdescribedusingSensitivity( **Sn**)andSpecificity( **Sp**)andthestatisticsforthegenelevelaregiven asMissedGene( **MG**)andWrongGene( **WG**).

# 9. References

Agarwal,P.andD.J.States.1998.Comparativeaccuracyofmethodsforproteinsequence
similaritysearch. *Bioinformatics* **14:**40-47.

Altschul,S.F.,W.Gish,W.Miller,E.W.MyersandD.J.Lipman.1990.Basiclocalalignment
searchtool. *JMolBiol* **215:**403-410.

Arkhipova,I.R.1995.PromoterelementsinDrosophilamelanogasterrevealedbysequenc      e
analysis. *Genetics* **139:**1359-1369.

Ashburner,M.2000. *submitted*.

Ashburner,M.ande.al.1999.EuropeanDrosophilaGenomeProject(EDGP).
http://edgp.ebi.ac.uk/.

Ashburner,M.,P.Bork,R.Durbin,R.GuigoandT.J.Hubbard.1999a.      *GASP1assessment
meeting,EMBL,Heidelberg* ,

Ashburner,M.,S.Misra,J.Roote,S.E.Lewis,R.Blazej,T.Davis,C.Doyle,R.Galle,R.
George,N.Harris,G.Hartzell,D.Harvey,L.Hong,K.
Houston,R.Hoskins,G.Johnson,C.Martin,A.Moshrefi,
M.Palazzolo,M.G.Reese,A.Spradling,G.Tsang,K.
Wan,K.Whitelaw,B.Kimmeland      *etal.* 1999b.An
explorationofthesequenceofa2.9-Mbregionofthe
genomeofdrosophilamelanogaster.Theadhregion.
*Genetics* **153:**179-219.

Bateman,A.,E.Birney,R.Durbin,S.R.Eddy,K.L.HoweandE.L.Sonnhammer.2000.The

PfamProteinFamiliesDatabase. *NucleicAcidsRes* **28:**
263-266.


Benson,G.1999.Tandemrepeatsfinder:aprogramtoanalyzeDNAsequences. *NucleicAcids
Res* **27:**573-580.


Besemer,J.andM.Borodovsky.1999.Heuristicapproachtoderivingmodelsforgenefinding.

*NucleicAcidsRes* **27:**3911-3920.


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


Birney,E.andR.Durbin.1997.Dynamite:aflexiblecodegeneratinglanguagefordynamic

programmingmethodsusedinsequencecomparison. *Ismb*
**5:**56-64.


Birney,E.andR.Durbin.2000.UsingGeneWiseintheDrosophilaannotationexperiment.

*GenomeResearch* **10**.


Burge,C.andS.Karlin.1997.PredictionofcompletegenestructuresinhumangenomicDNA. *J
MolBiol* **268:**78-94.


Burge,C.B.andS.Karlin.1998.FindingthegenesingenomicDNA. *CurrOpinStructBiol* **8:**
346-354.


Burset,M.andR.Guigo.1996.Evaluationofgenestructurepredictionprograms. *Genomics* **34:**
353-367.

CavinPerier,R.,T.Junier,C.BonnardandP.Bucher.1999.TheEukaryoticPromoterDatabase
(EPD):recentdevelopments. *NucleicAcidsRes* **27:**307-309.

CavinPérier,R.,V.Praz,T.Junier,C.BonnardandP.Bucher.2000.TheEukaryoticPromoter
Database(EPD). *NucleicAcidsRes* **28:**302-303.

Dunbrack,R.L.,Jr.,D.L.Gerloff,M.Bower,X.Chen,O.LichtargeandF.E.Cohen.1997.
Meetingreview:theSecondmeetingontheCritical
AssessmentofTechniquesforProteinStructurePrediction
(CASP2),Asilomar,California,December13-16,1996.
*FoldDes* **2:**R27-42.

Eeckman,F.H.andR.Durbin.1995.ACeDBandmacace. *MethodsCellBiol* **48:**583-605.

Fickett,J.W.andA.G.Hatzigeorgiou.1997.Eukaryoticpromoterrecognition. *GenomeRes* **7:**861-878.

Fickett,J.W.andC.S.Tung.1992.Assessmentofproteincodingmeasures. *NucleicAcidsRes*
**20:**6441-6450.

Field,D.1999. *unpublished.*

Florea,L.,G.Hartzell,Z.Zhang,G.M.RubinandW.Miller.1998.Acomputerprogramfor
aligningacDNAsequencewithagenomicDNA
sequence. *GenomeRes* **8:**967-974.

Friese,E.,M.G.ReeseandG.M.Rubin.1999.      *ProceedingsoftheThirdAnnualInternational ConferenceonComputationalMolecularBiology (RECOMB),Lyon,France* ,

Gaasterland,T.andC.W.Sensen.1996.MAGPIE:automatedgenomeinterpretation.      *Trends Genet* **12:**76-78.

Green,P.1995.   *unpublished.*

Guigo,R.,S.Knudsen,N.DrakeandT.Smith.1992.Predictionofgenestructure.      *JMolBiol* **226:**141-157.

Harris,N.L.,G.Helt,S.MisraandS.E.Lewis.1999.CloneCurator. http://www.fruitfly.org/displays/CloneCurator.html.

Helt,G.ande.al.1999.NeomorphicGenomeSoftwareDevelopmentToolkit(NGSDK). NeomorphicInc.,Berkeley.http://www.neomorphic.com.

Henikoff,J.G.,S.HenikoffandS.Pietrokovski.1999a.NewfeaturesoftheBlocksDataba        se servers. *NucleicAcidsRes*   **27:**226-228.

Henikoff,S.andJ.G.Henikoff.1994a.Proteinfamilyclassificationbasedonsearchinga databaseofblocks.   *Genomics* **19:**97-107.

Henikoff,S.andJ.G.Henikoff.1994b.      *27thAnn.HawaiiIntl.ConferenceonSystemSciences, Hawaii,U.S.A.* ,

Henikoff,S.andJ.G.Henikoff.2000.Genomicsequenceannotationbasedontranslated

searchingoftheBlocks+Database. *GenomeResearch* .

Henikoff,S.,J.G.HenikoffandS.Pietrokovski.1999b.Blocks+:anon-redundantdatabaseof

proteinalignmentblocksderivedfrommultiple

compilations. *Bioinformatics* **15:**471-479.

Hubbard,T.J.2000.Personalcommunication..

Jurka,J.1998.RepeatsingenomicDNA:miningandmeaning. *CurrOpinStructBiol* **8:**333-

337.

Krogh,A.1997.TwomethodsforimprovingperformanceofanHMMandtheirapplicationfor

genefinding. *Ismb* **5:**179-186.

Kulp,D.,D.Haussler,M.G.ReeseandF.H.Eeckman.1997.Integratingdatabasehomology ina

probabilisticgenestructuremodel. *PacSympBiocomput* **:**

232-244.

Kurtz,S.andC.Schleiermacher.1999.REPuter:fastcomputationofmaximalrepeatsin

completegenomes. *Bioinformatics* **15:**426-427.

Levitt,M.1997.Competitiveassessmentofproteinfoldrecognitionandalignmentaccur acy.

*Proteins* **Suppl:**92-104.

Marcotte,E.M.,M.Pellegrini,M.J.Thompson,T.O.YeatesandD.Eisenberg.1999.A

combinedalgorithmforgenome-widepredictionofprotein

function. *Nature* **402:**83-86.

Mott,R.1997.EST_GENOME:aprogramtoalignsplicedDNAsequencestounsplicedgenom ic
DNA. *ComputApplBiosci* **13:**477-478.

Moult,J.,T.Hubbard,S.H.Bryant,K.FidelisandJ.T.Pedersen.1997.Criticalassessmentof
methodsofproteinstructureprediction(CASP):roundII.
*Proteins* **Suppl:**2-6.

Moult,J.,T.Hubbard,K.FidelisandJ.T.Pedersen.1999.Criticalassessmentofmethodsof
proteinstructureprediction(CASP):roundIII. *Proteins*
**Suppl:**2-6.

Ohler,U.,S.Harbeck,H.Niemann,E.NothandM.G.Reese.1999.Interpolatedmarkovchains
foreukaryoticpromoterrecognition. *Bioinformatics* **15:**
362-369.

Ohler,U.,G.StommerandS.Harbeck.2000.StochasticSegmentModelsofEukaroyotic
PromoterRegions. *PacSympBiocomput* **5:**377-388.

Parra,G.,E.BlancoandR.Guigo.2000.GeneIDinDrosophila. *GenomeResearch* **10**.

Pearson,W.R.1995.Comparisonofmethodsforsearchingproteinsequencedatabases. *Protein
Sci* **4:**1145-1160.

Pearson,W.R.andD.J.Lipman.1988.Improvedtoolsforbiologicalsequencecomparison. *Proc
NatlAcadSciUSA* **85:**2444-2448.

Reese,M.G.2000.GenomeAnnotationin *Drosophilamelanogaster* .Ph.D.,Universityof
Hohenheim.

Reese,M.G.,F.H.Eeckman,D.KulpandD.Haussler.1997.Improvedsplicesitedetectionin

Genie. *JComputBiol* **4:**311-323.

Reese,M.G.,N.L.Harris,G.HartzellandS.E.Lewis.1999.                 *The7thconferenceonIntelligent*

*SystemsinMolecularBiology(ISMB'99),Heidelberg,*

*Germany*,http://www.fruitfly.org/GASP1.

Reese,M.G.,D.Kulp,H.TammanaandD.Haussler.2000.Genie-Genefindingin                 *Drosophila*

*melanogaster*. *GenomeResearch* **10**.

Rubin,G.M.2000.Full-lengthcDNAproject..

Rubin,G.M.ande.al.1999.BerkeleyDrosophiaGenomeProject(BDGP).

http://www.fruitfly.org.

Salamov,A.A.andV.V.Solovyev.2000.AbinitiogenefindinginDrosophilagenomicDNA.

*GenomeResearch* **10**.

Scherf,M.,A.KlingenhoffandT.Werner.2000.         *inpreparation.*

Sippl,M.J.,P.Lackner,F.S.DominguesandW.A.Koppensteiner.1999.Anattempttoanalyse

progressinfoldrecognitionfromCASP1toCASP3.

*Proteins* **Suppl:**226-230.

Solovyev,V.V.,A.A.SalamovandC.B.Lawrence.1995.Identificationofhumangene

structureusinglineardiscriminantfunctionsanddynamic

programming. *Ismb* **3:**367-375.

Sonnhammer,E.L.,S.R.Eddy,E.Birney,A.BatemanandR.Durbin.1998.Pfam:multiple

sequencealignmentsandHMM-profilesofprotein

domains. *NucleicAcidsRes* **26:**320-322.

Sonnhammer,E.L.,S.R.EddyandR.Durbin.1997.Pfam:acomprehensivedatabaseofprotein

domainfamiliesbasedonseedalignments. *Proteins* **28:** 405-420.

Stein,L.D.andJ.Thierry-Mieg.1998.ScriptableaccesstotheCaenorhabditiselegans genome

sequenceandotherACEDBdatabases. *GenomeRes* **8:** 1308-1315.

Stormo,G.D.2000. *submitted*.

Uberbacher,E.C.andR.J.Mural.1991.Locatingprotein-codingregionsinhumanDNA

sequencesbyamultiplesensor-neuralnetworkapproach.

*ProcNatlAcadSciUSA* **88:**11261-11265.

Zemla,A.,C.Venclovas,J.MoultandK.Fidelis.1999.ProcessingandanalysisofCASP3

proteinstructurepredictions. *Proteins* **Suppl:**22-29.