



15K WEBINAR:
AUTOMATED GENOME
ANNOTATION AND ANALYSIS

Carson Holt

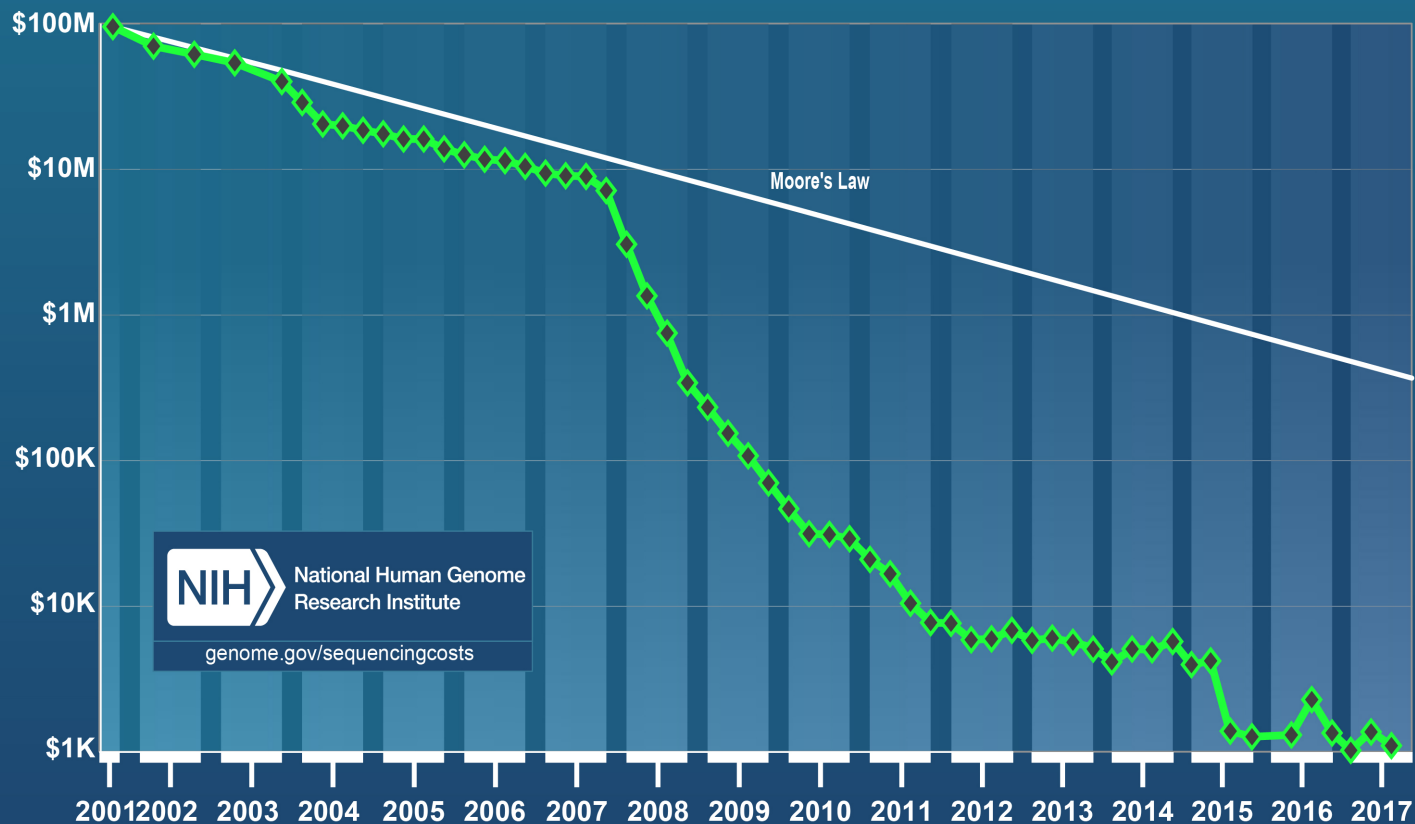
Yandell Lab

USTAR Center for Genetic Discovery

University of Utah

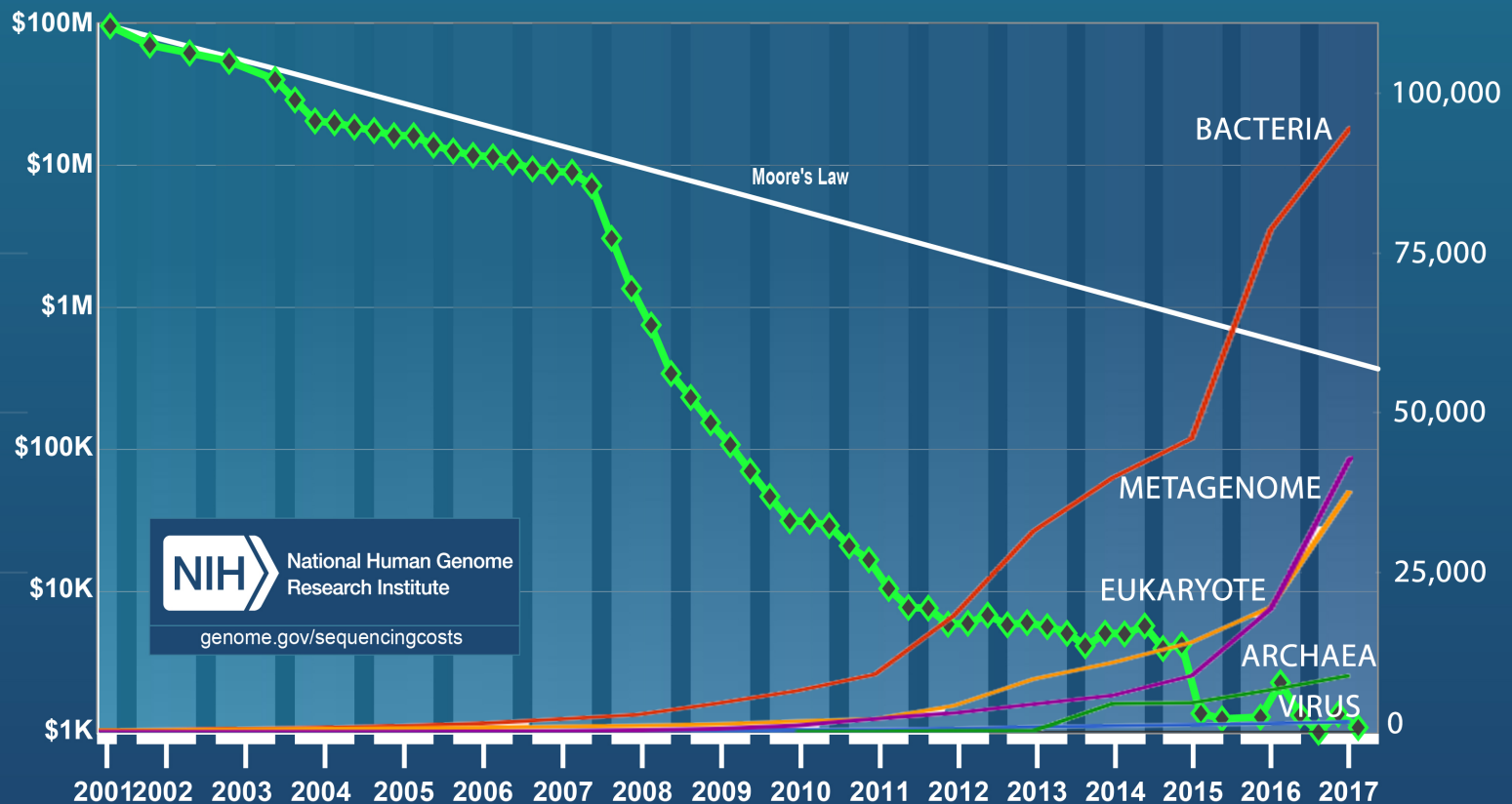
Advances in Second-Generation Technology are Making Whole Genome and Transcriptome Sequencing “Routine” Even for Small Labs

Cost per Genome



Advances in Second-Generation Technology are Making Whole Genome and Transcriptome Sequencing “Routine” Even for Small Labs

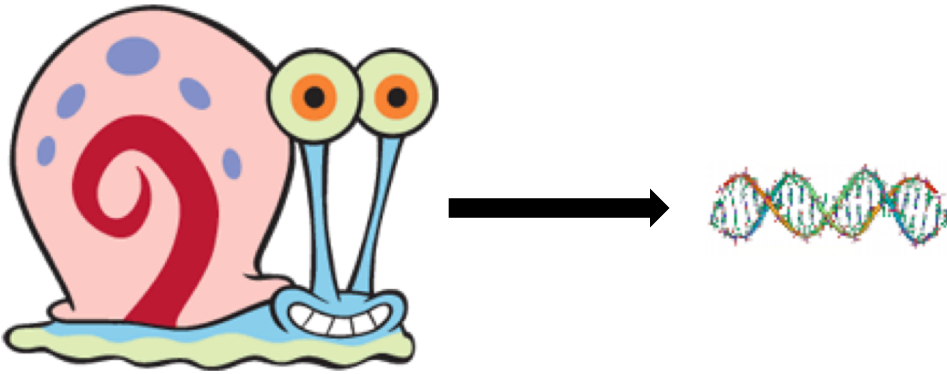
Cost per Genome



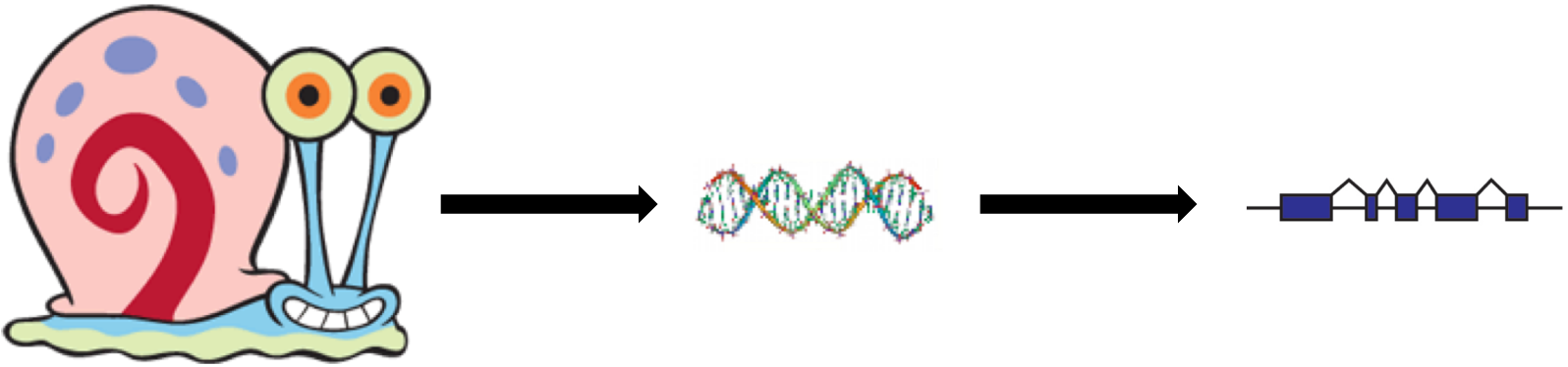
Genome Project Overview



Genome Project Overview

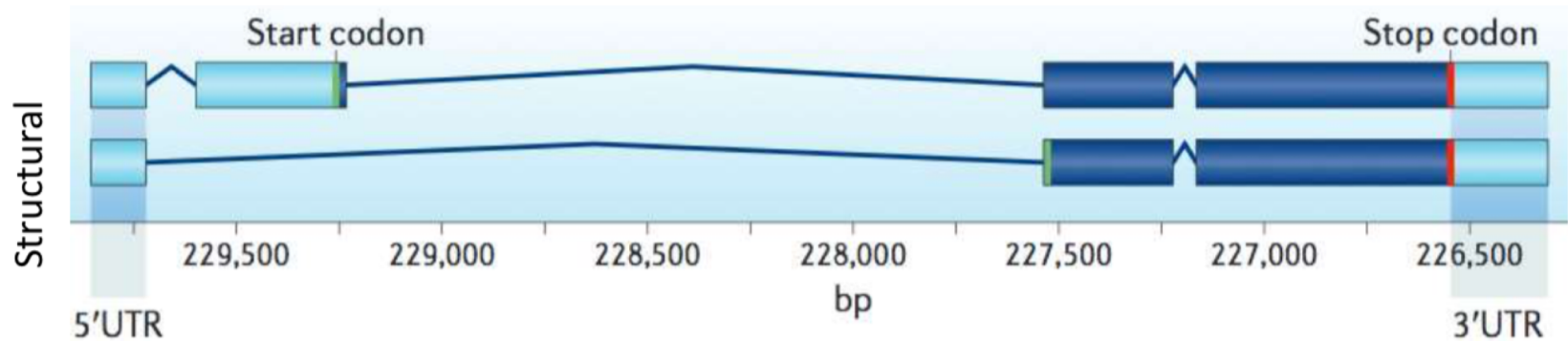


Genome Project Overview



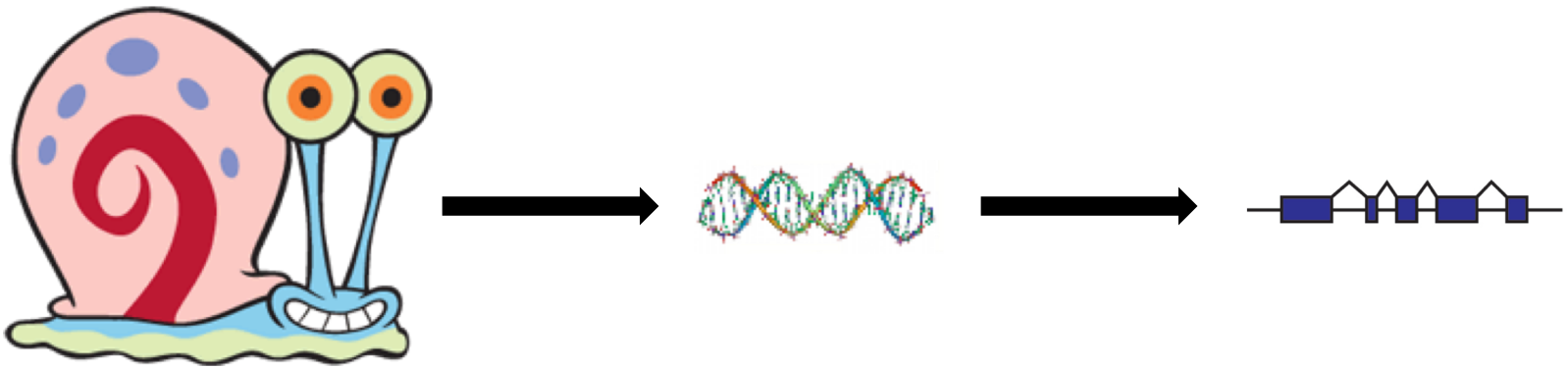
Genome Project Overview

What is an Annotation?

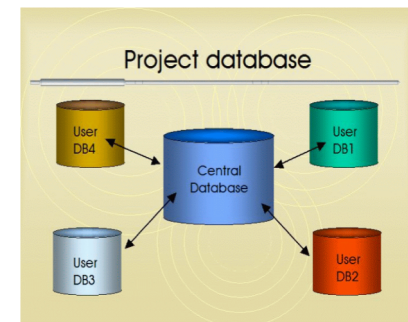
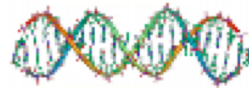


Functional	Function	cAMP-dependent and sulfonyleurea-sensitive anion transporter. Key gatekeeper influencing intracellular cholesterol transport.
	Subcellular location	Membrane; Multi-pass membrane protein Ref.13 Ref.14.
	Domain	Multifunctional polypeptide with two homologous halves, each containing a hydrophobic membrane-anchoring domain and an ATP binding cassette (ABC) domain.

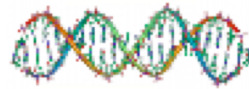
Genome Project Overview



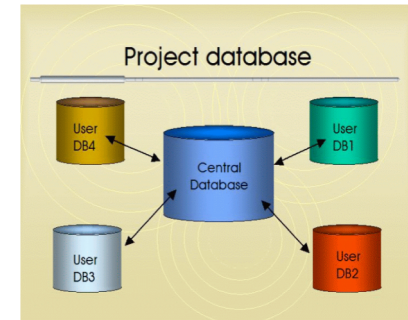
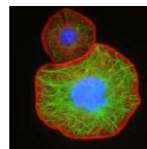
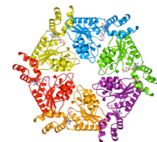
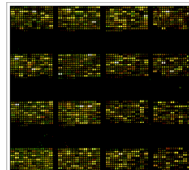
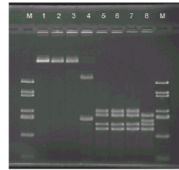
Genome Project Overview



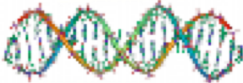
Genome Project Overview



```
>Smg5  
MEVTFSSGGSSNASSECAIDGGTNRCRGL  
EPNNGTCILSQEVKDLYRSLYTASKQLDD  
AKRNVQSVGQLFQHEIEEKRSLLVQLCKQ  
IIFKDYQSVGKKVREVMWRRGYEFIAFV
```

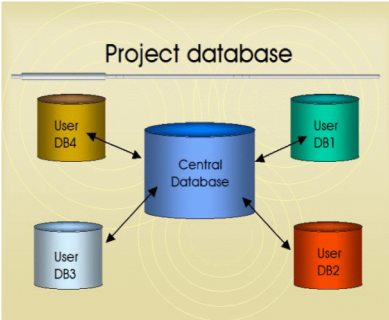
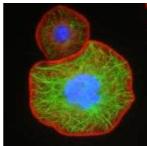
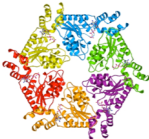
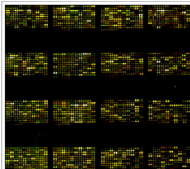
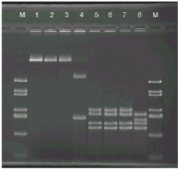


Genome Project Overview



```
>Smg5  
MEVTFSSGGSSNASSECAIDGGTNRCRGL  
EPNNGTCILSQEVKDLYRSLYTASKQLDD  
AKRNVQSVGQLFQHEIEEKRSLLVQLCKQ  
IIFKDYQSVGKKVREVMWRRGYEFIAFV
```

SUCCESS



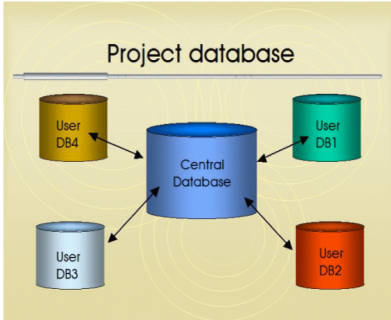
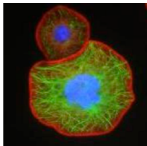
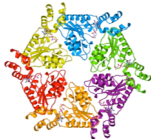
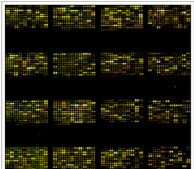
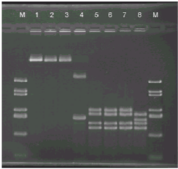
Genome Project Overview



Incorrect annotations poison every experiment that uses them!!

```
>Smg5  
MEVTFSSGGSSNASSECAIDGGTNRCRGL  
EPNNGTCILSQEVKDLYRSLYTASKQLDD  
AKRNVQSVGQLFQHEIEEKRSLLVQLCKQ  
IIFKDYQSVGKKVREVMWRRGYEFIAFV
```

FAILURE





MAKER

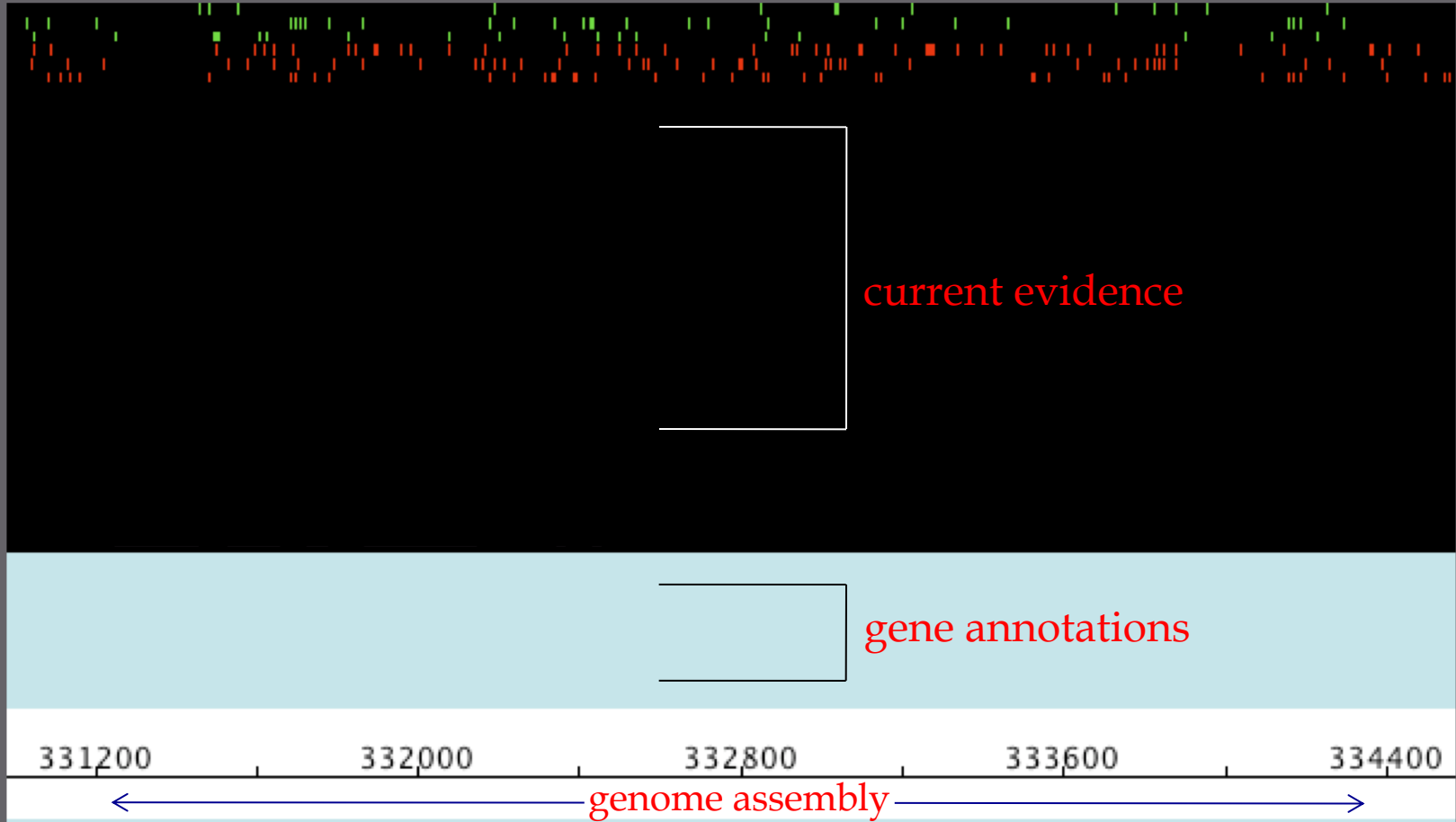
an annotation pipeline and genome-database management tool for second-generation genome projects

Easy-to-use by design

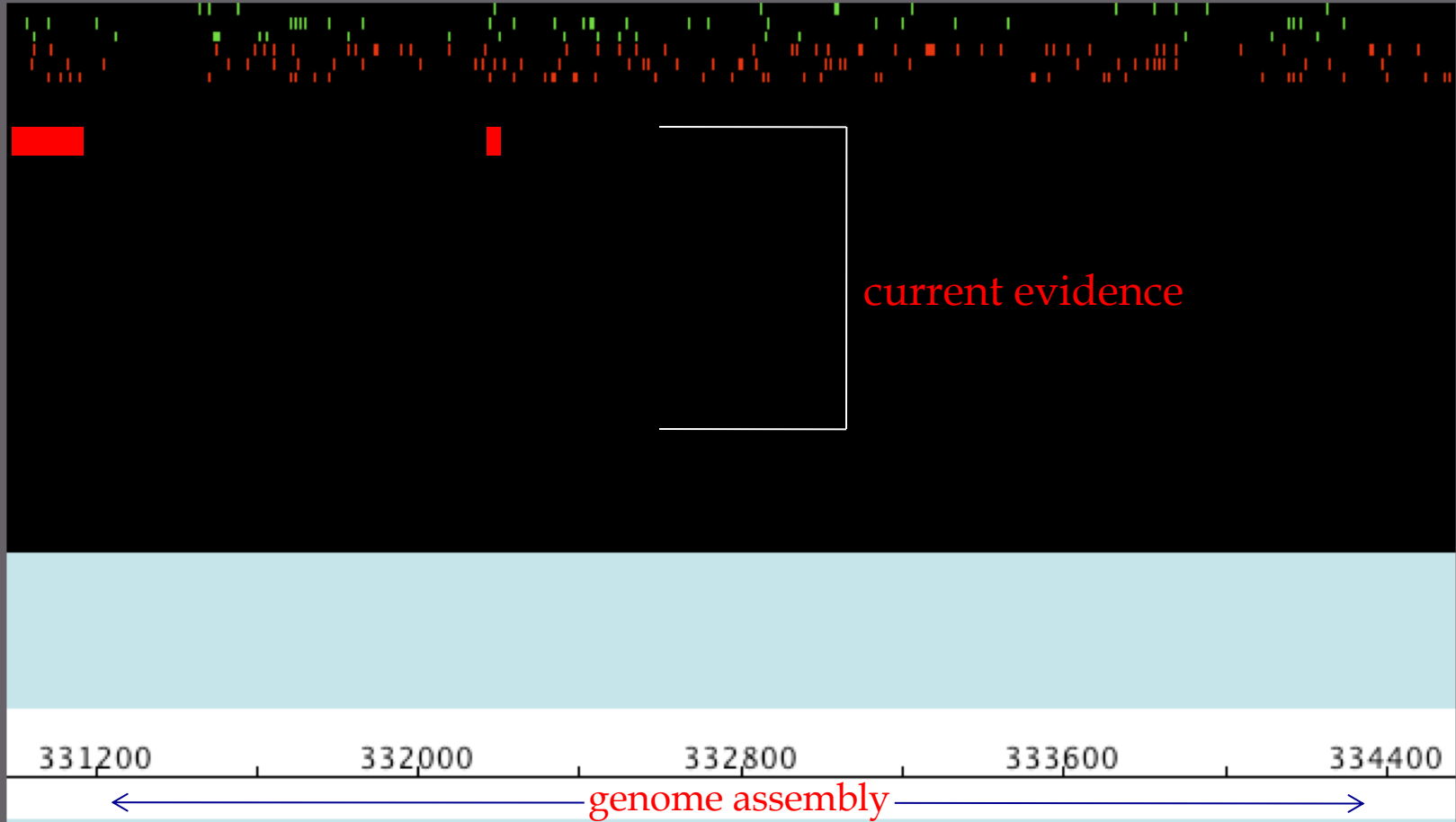
User Requirements:	Can be run by a single individual with little bioinformatics experience
System Requirements:	Can run on laptop or desktop computers (running Linux or Mac OS X)
Program Output:	Output is compatible with popular annotation tools like Apollo, GBrowse, and JBrowse
Availability:	Free open source application (for the academic community)

How does MAKER work?

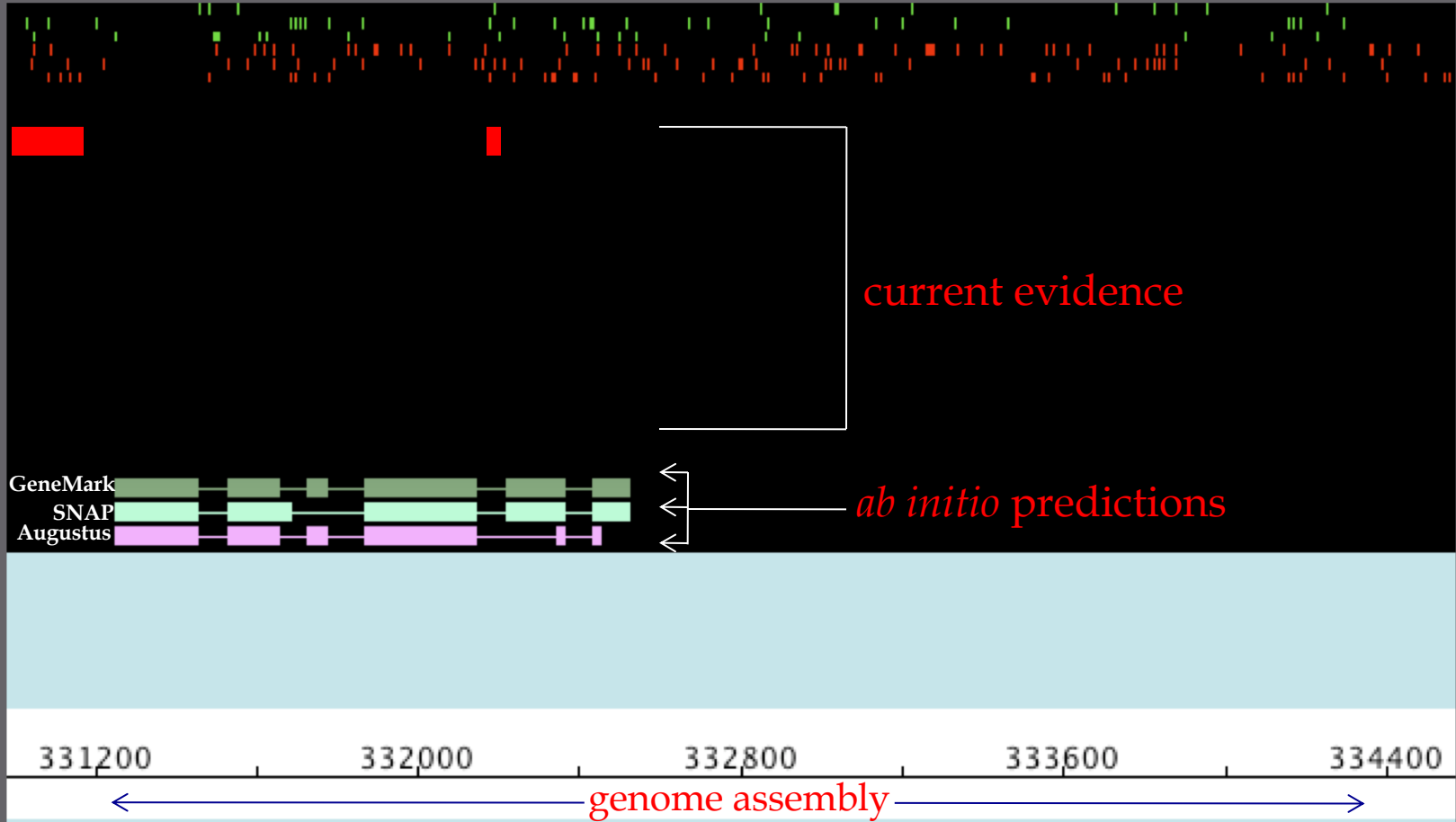
Annotating the Genome – Apollo View



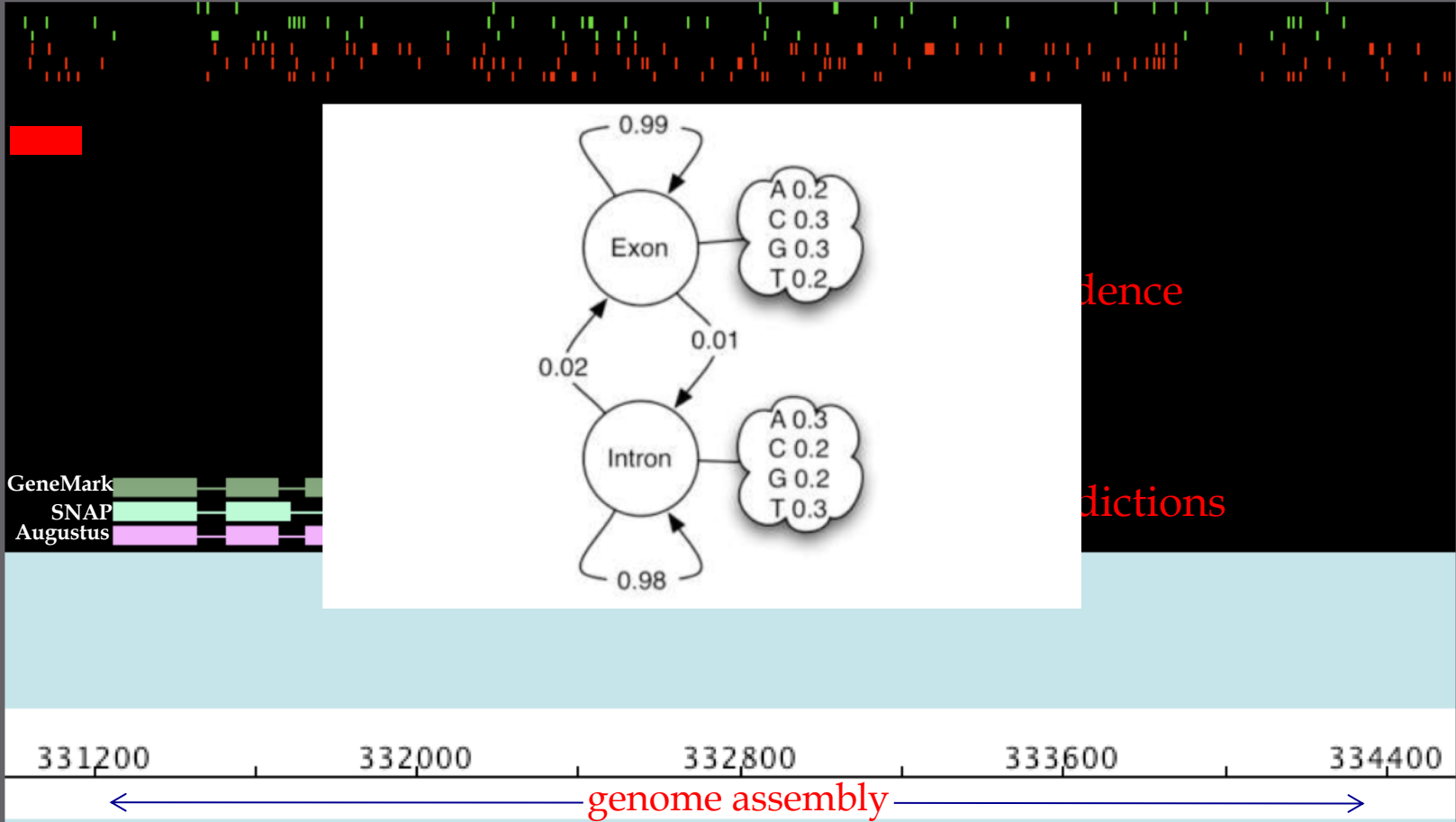
Identify and mask repetitive elements



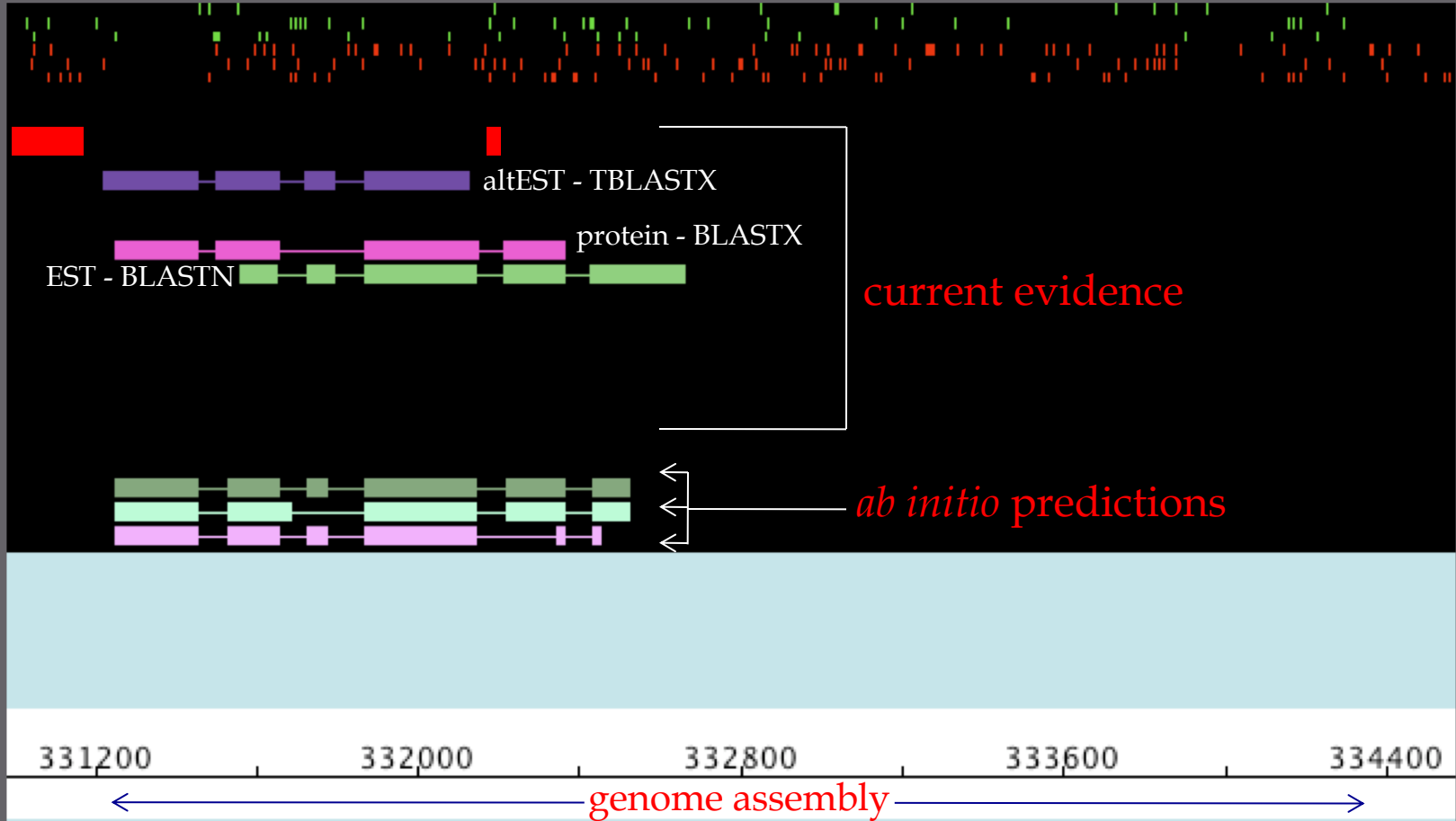
Generate *ab initio* gene predictions



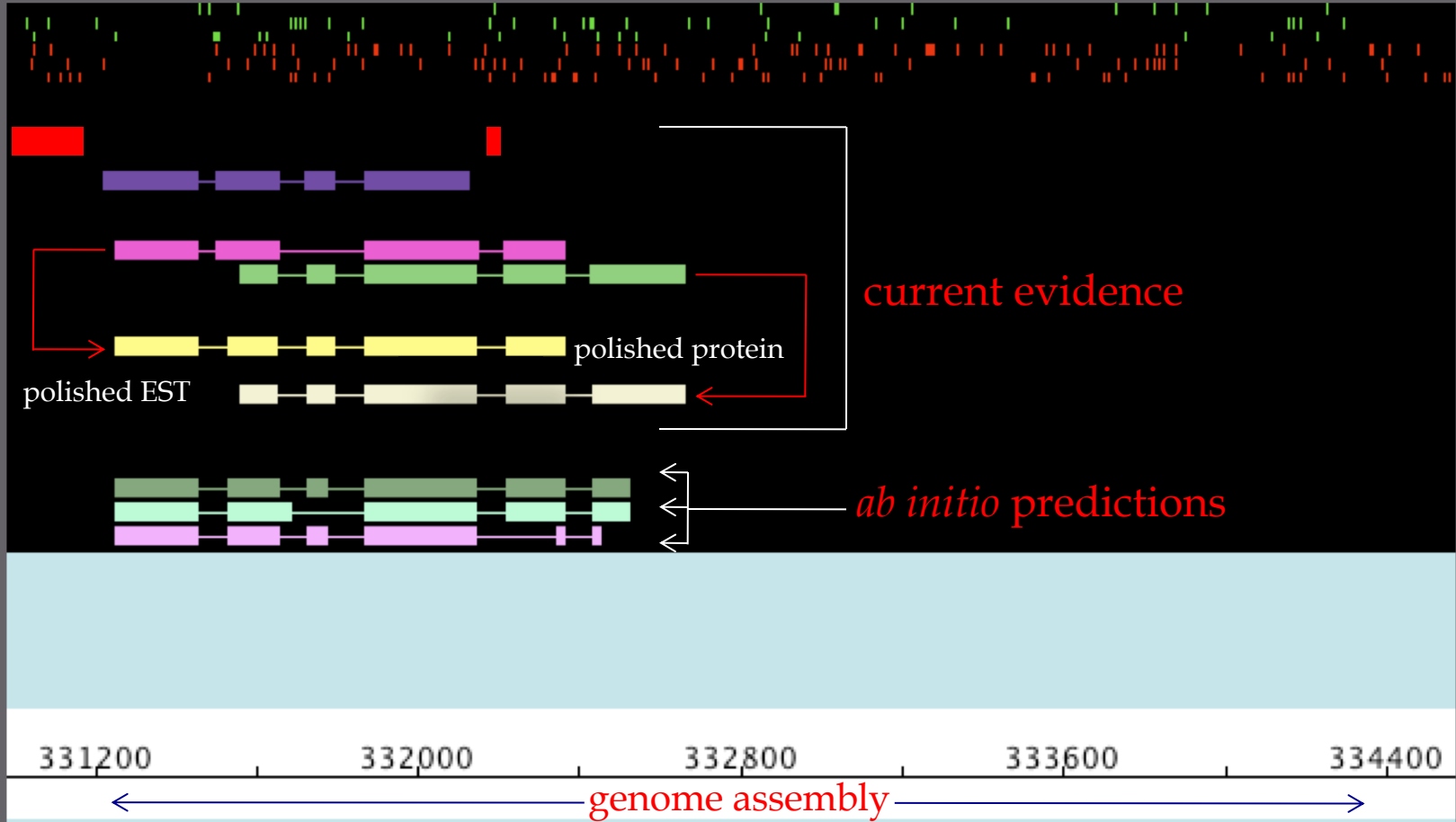
Generate *ab initio* gene predictions



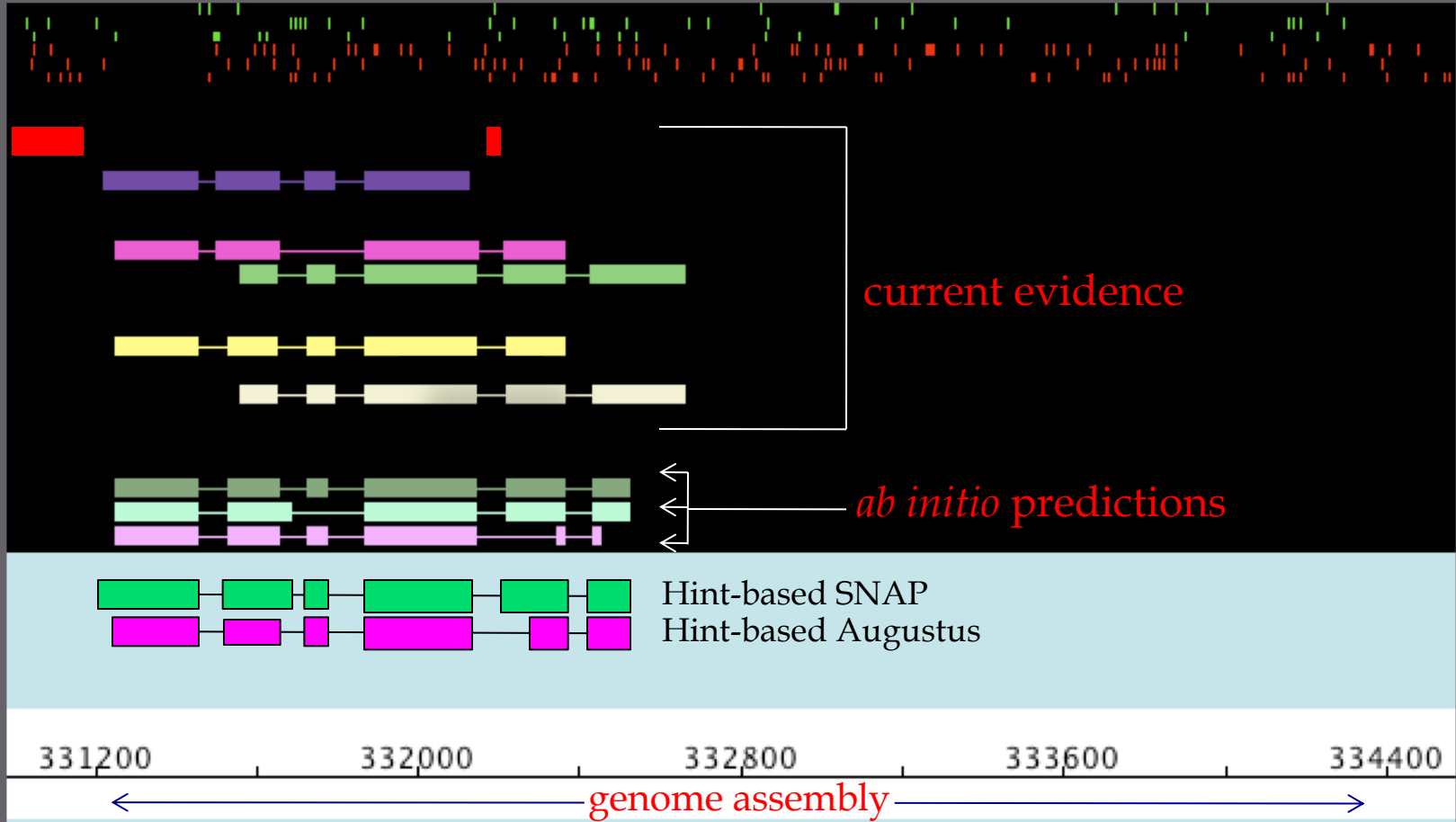
Align EST and protein evidence



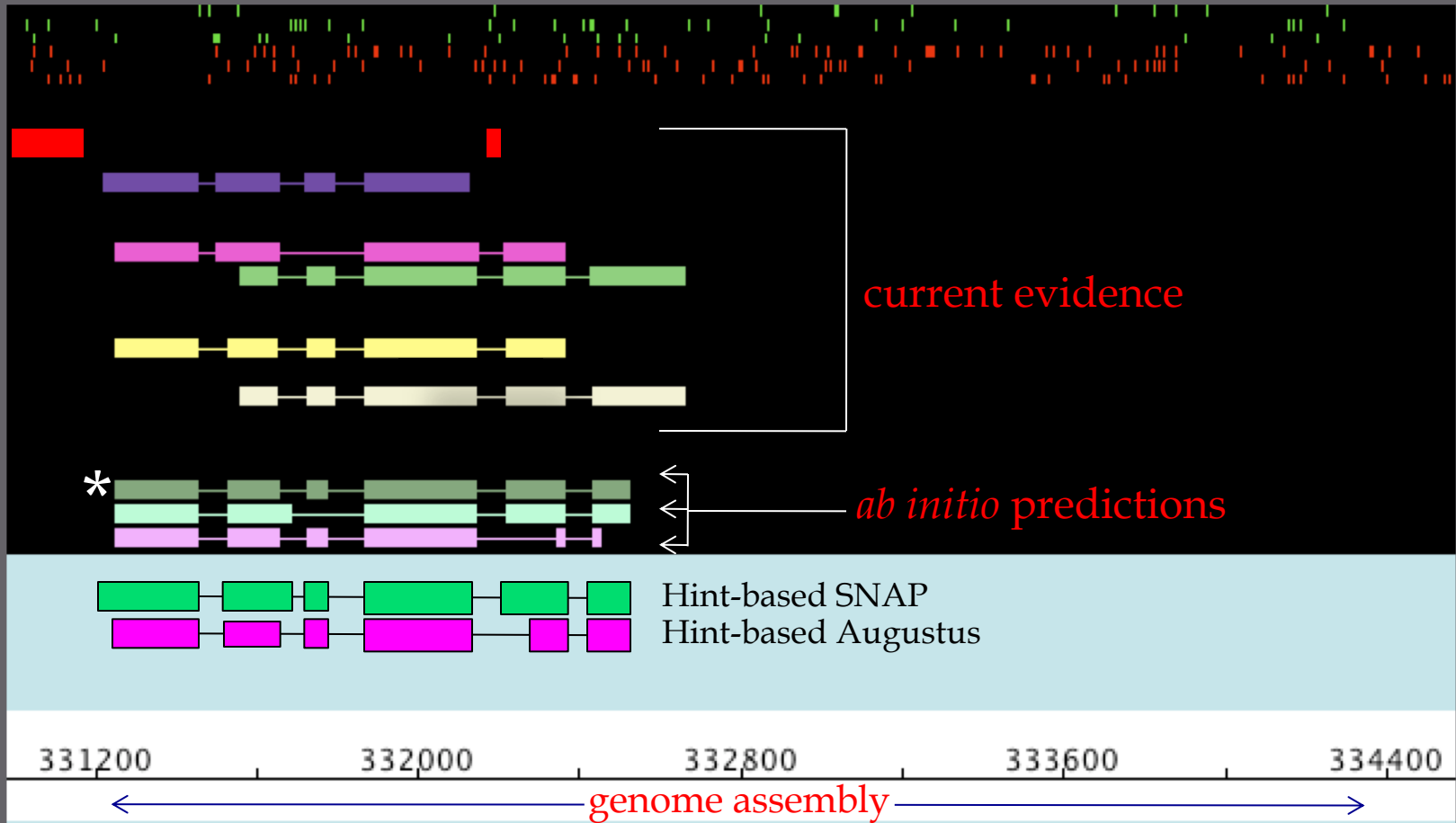
Polish BLAST alignments with Exonerate



Pass gene-finders evidence-based 'hints'

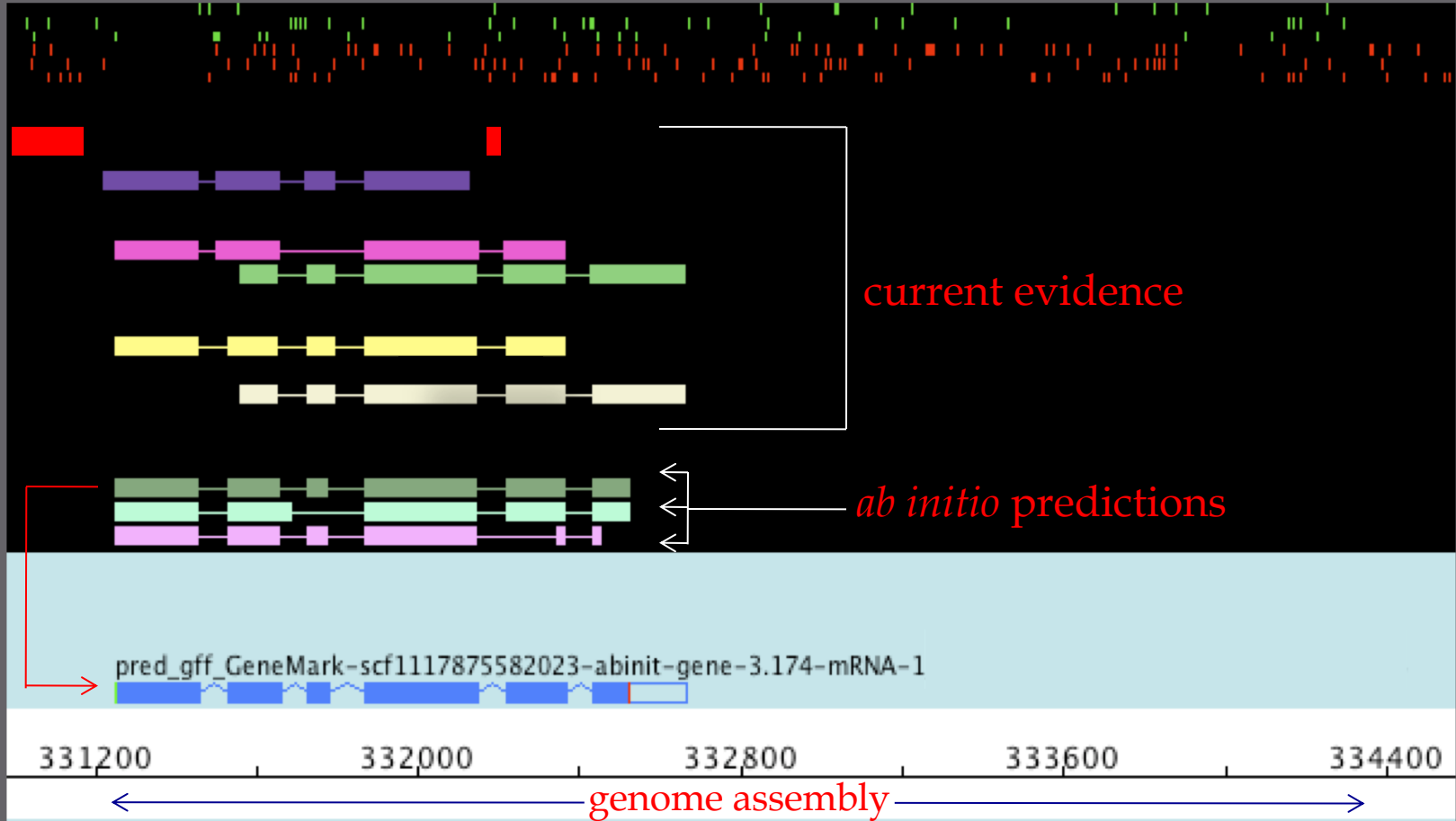


Identify gene model most consistent with evidence

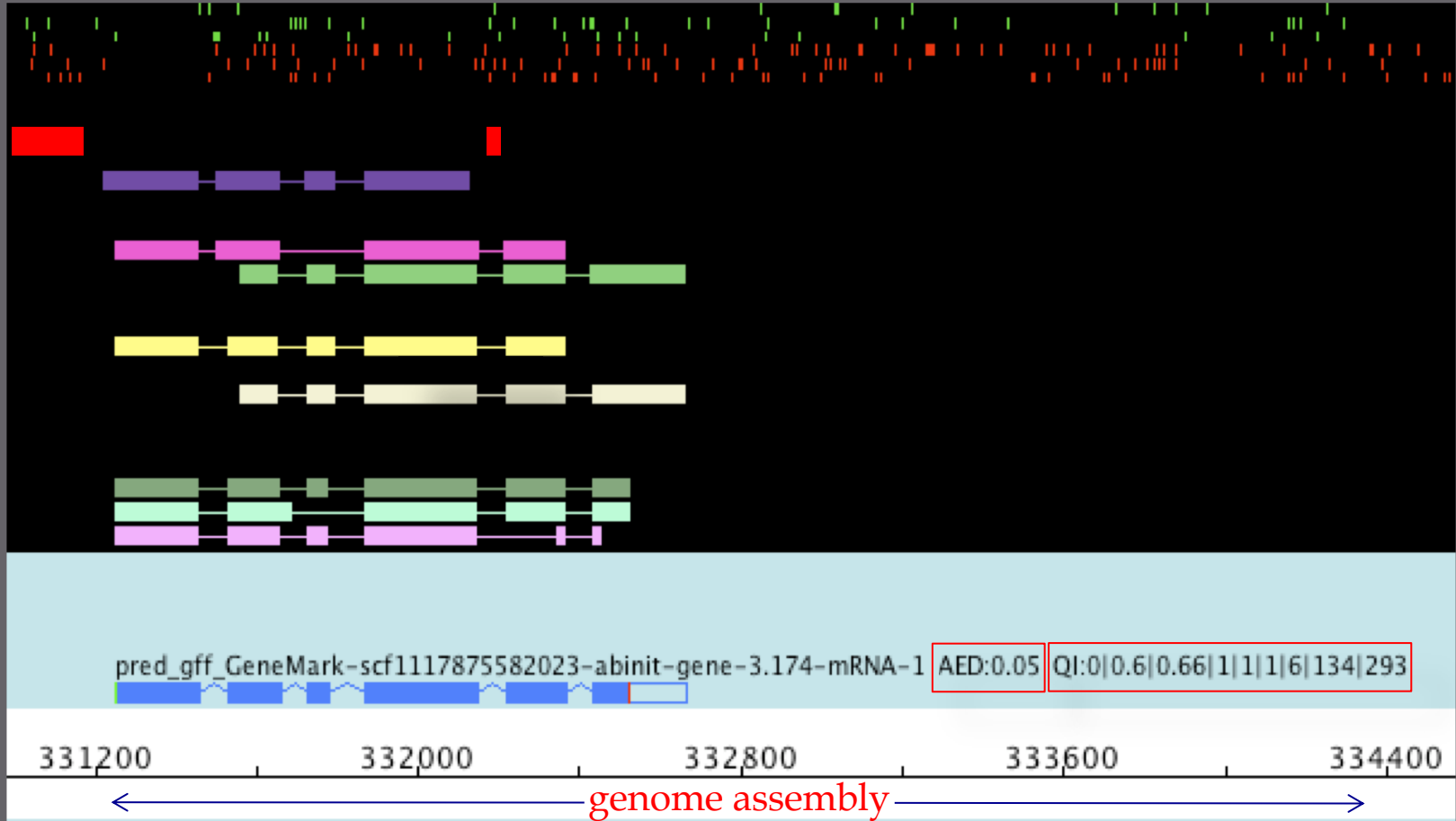


*Quantitative Measures for the Management and Comparison of Annotated Genomes
Karen Eilbeck , Barry Moore , Carson Holt and Mark Yandell BMC Bioinformatics 2009
10:67doi:10.1186/1471-2105-10-67

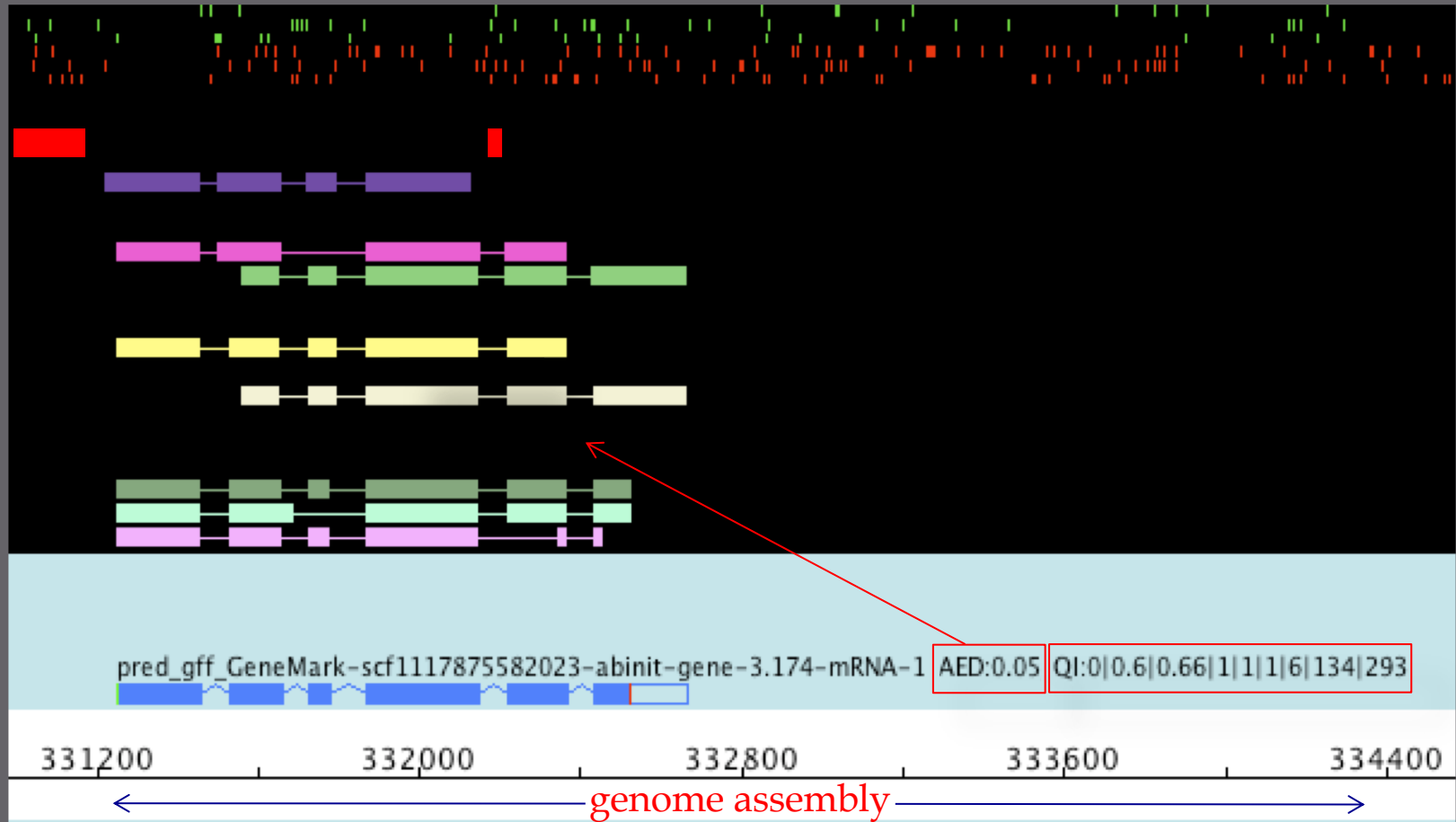
Revise it further if necessary; create new annotation



Compute support for each portion of gene model



Compute support for each portion of gene model



*Quantitative Measures for the Management and Comparison of Annotated Genomes
Karen Eilbeck , Barry Moore , Carson Holt and Mark Yandell BMC Bioinformatics 2009
10:67doi:10.1186/1471-2105-10-67

Compute support for each portion of gene model



*Cantarel BL, Korf I, Robb SMC, Parra G, Ross E, Moore B, Holt C, Sanchez Alvarado A, Yandell M: **MAKER: An easy-to-use annotation pipeline designed for emerging model organism genomes.** *Genome Res* 2008, 18:188-196.

```
bin -- bmoore@derringer:/data1/genomes/Homo/sapiens/36.2 -- ssh
```

scaffold00080	repeatmasker	match_part	723561	723736	914	+	.	ID=scaffold00080:hsp:4987;Par
scaffold00080	repeatmasker	match	724950	725171	495	+	.	ID=scaffold00080:hit:3405;Name=specie
scaffold00080	repeatmasker	match_part	724950	725171	495	+	.	ID=scaffold00080:hsp:4988;Par
scaffold00080	repeatmasker	match	726925	727293	724	+	.	ID=scaffold00080:hit:3406;Name=specie
scaffold00080	repeatmasker	match_part	726925	727293	724	+	.	ID=scaffold00080:hsp:4989;Par
scaffold00080	maker	gene	56197	58302	.	+	.	ID=ACEP_00015614;Name=ACEP_00015614;Alias=mak
scaffold00080	maker	mRNA	56197	58302	.	+	.	ID=ACEP_00015614-RA;Parent=ACEP_00015614;Name
scaffold00080	maker	exon	56197	56274	.	+	.	ID=ACEP_00015614-RA:exon:126;Parent=ACEP_0001
scaffold00080	maker	exon	56569	56584	.	+	.	ID=ACEP_00015614-RA:exon:127;Parent=ACEP_0001
scaffold00080	maker	exon	56797	56906	.	+	.	ID=ACEP_00015614-RA:exon:128;Parent=ACEP_0001
scaffold00080	maker	exon	57851	57941	.	+	.	ID=ACEP_00015614-RA:exon:129;Parent=ACEP_0001
scaffold00080	maker	exon	58067	58302	18.726	+	.	ID=ACEP_00015614-RA:exon:130;Parent=ACEP_0001
scaffold00080	maker	CDS	56197	56274	.	+	0	ID=ACEP_00015614-RA:cds:124;Parent=ACEP_00015
scaffold00080	maker	CDS	56569	56584	.	+	0	ID=ACEP_00015614-RA:cds:125;Parent=ACEP_00015
scaffold00080	maker	CDS	56797	56906	.	+	2	ID=ACEP_00015614-RA:cds:126;Parent=ACEP_00015
scaffold00080	maker	CDS	57851	57941	.	+	1	ID=ACEP_00015614-RA:cds:127;Parent=ACEP_00015
scaffold00080	maker	CDS	58067	58302	.	+	2	ID=ACEP_00015614-RA:cds:128;Parent=ACEP_00015
scaffold00080	maker	gene	2797	3250	.	-	.	ID=ACEP_00015615;Name=ACEP_00015615;Alias=mak

GFF3

```
bin -- bmoore@derringer:/data1/genomes/Homo/sapiens/36.2 -- ssh
```

```
>ACEP_00015614-RA protein AED:0.401129943502825 QI:0|0|0|0.2|1|1|5|0|177
MEKSDDYQDHYVIQFLVYMYFVHEKQYYLWLVKRIITLSNKYYPIGVIWAFFVPIVF
IRHPDDLKTIILSNPKHIKKSFFYDNIKPWLGTSILTSEDCTLSHFIPMLSNRHFTIVHCI
VHVLGTKWQLQRKILISTFHFIDILNQFVEIFEKESKNMIKSLKNAEGTVVKDLSSFI
>ACEP_00015615-RA protein AED:0.223684210526316 QI:60|0|0|0.5|0|0|2|0|93
MDTQQKHKREIPTPEGEINSVISFQICISCCDLSYCNIESPTNATNATYISRRRAKSKSK
RKRPRGRNGVDAATRLYGSSLLWLPASLFYAYHC
>ACEP_00015613-RA protein AED:0.191135958515638 QI:9|0.4|0.33|0.83|0.2|0.5|6|0|238
MEKSDNSLTYQDHYISIFGSVYVLFGRSREILWKIITLSNKFLRTFYPIGIIIRAFFVPIVSI
RHPDDLKSFYDNIKPWLGTSILINEGTKWQLQRKILIPTFHFIDILNQFVEIFEKESKNM
IKSLKNAEGTVVKDLSSFISEYTLNAICESYLIPAGTVLHININGVHTDPNFWRYFQIQR
YLILIDFCPRGSEIVTLTHIYHSVLDHVIVSNLRFKSKFDHERSGIVREI
>ACEP_00015612-RA protein AED:0.469714351691491 QI:249
MYFVNFNYFYANASVTNHRNGRTSLHTELLSSKVKREQTLRNILTARRE
RILSSRRRPRPGDIISKDPQDVLETQETRNGIYIANARTEISIKEEIPQ
INLNATGEELESRTTKSKNESLSPTTEKIENISSIVSTPATNQKLGDAEMRKKNTVSMF
ESLYNHFRPVESNIPVEDMSQFLYFGQKLQPDALNVTSSSNNSVETTSTNPTSSRRRYST
KRFTATIATPMEIIMNEEMNIALETLEERRTVEKNQVSRKNTFRSSGNGLYYRKRPTAD
VVSNIIPGNESSPGIIGGSETKSEMNTVDDKFHRSDSSAHADGSRDRVPLERKNAHVSE
VSMTTQIPRQSDVEARVIVESEIGAKLSNKTTITASNGIVESLQRINDFGKDNEKKKITEA
EAGNVEKSEHKLMSVERTISTSQRESERFRESDESTMSLVPSRAESSTLRIVVPDNSFG
LCCFATEGECRCHUNTYKQDVAWDDTADRDERRRCCCAACCCGAEITTCGCCGCRDQCCQDRD
```

FASTA



GMOD in the Cloud toolset



GBrowse: Genome annotation viewer



Galaxy: Data analysis & integration



Chado: Biological database schema



JBrowse: Super-fast genome annotation viewer



BioMart: Data mining system



WebApollo: browser-based annotation editor



MAKER: Genome annotation pipeline



GBrowse_syn: Synteny viewer



Tripal: Chado web interface



InterMine: Data warehousing



CMap: Comparative map viewer



Pathway Tools: Metabolic, regulatory pathways

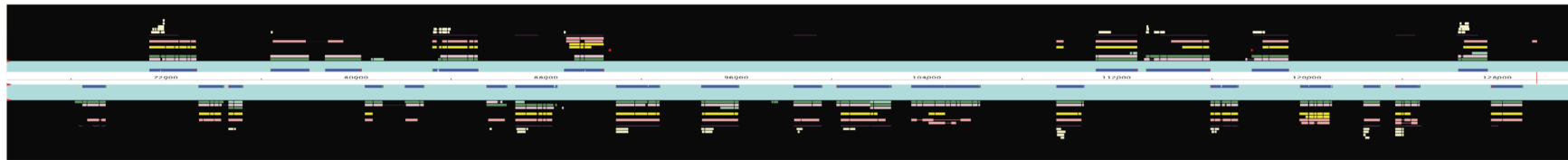
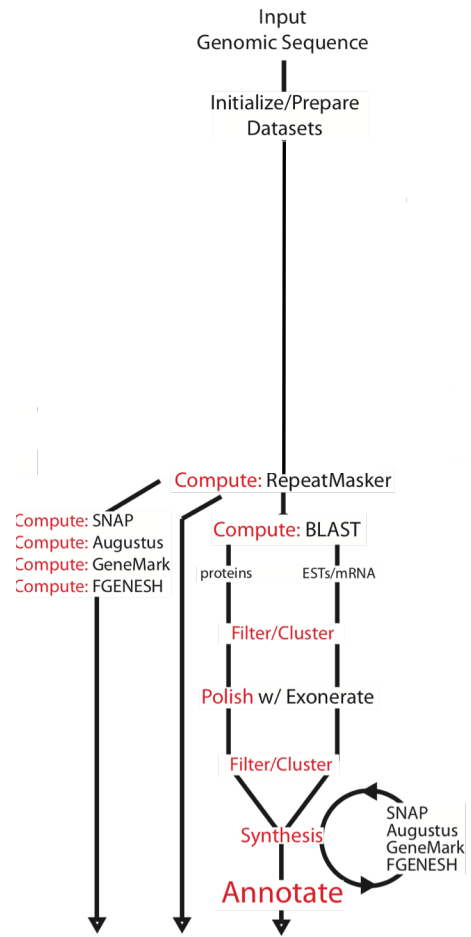


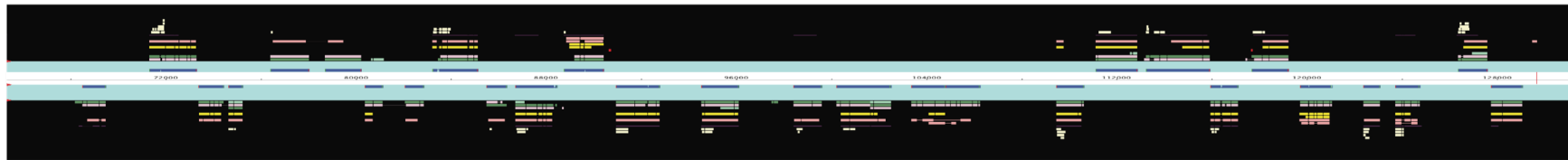
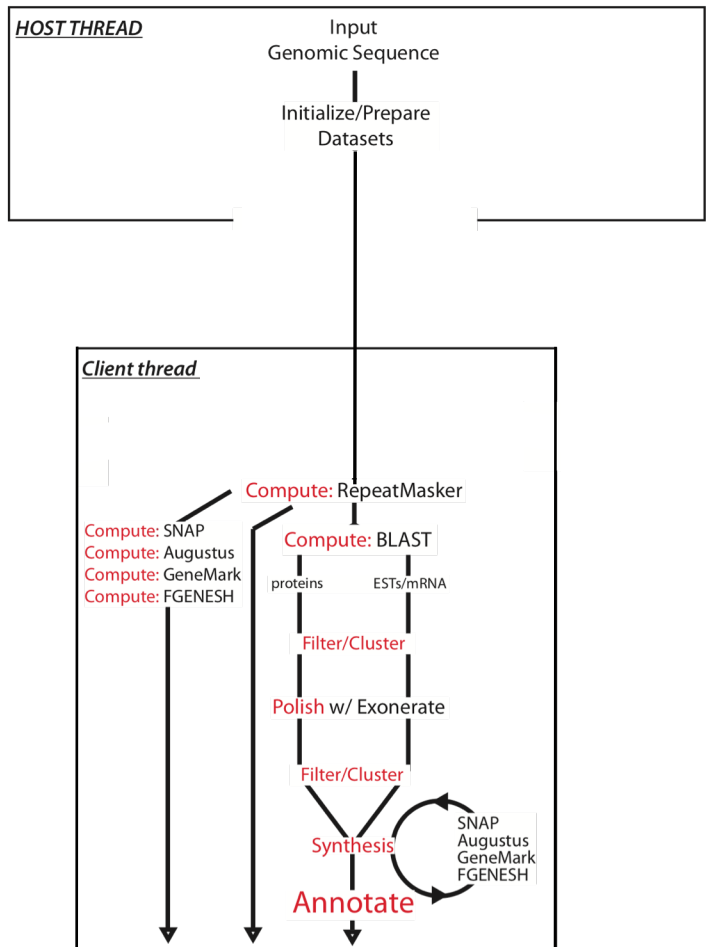
Canto: literature annotation tool

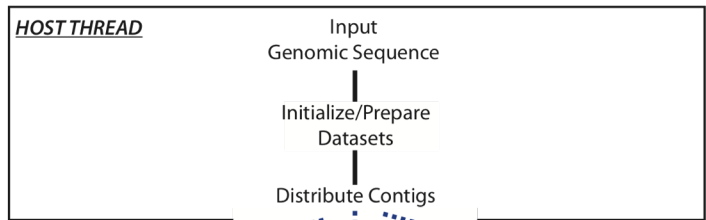
Distributed Parallelization

- Supports Message Passing Interface (MPI), a communication protocol for computer clusters which essentially allows multiple computers to act like a single powerful machine.

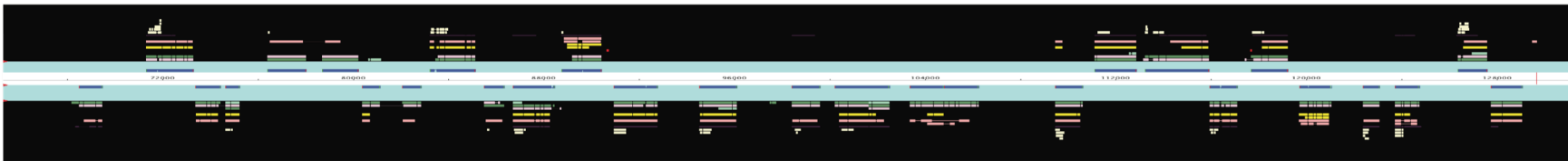
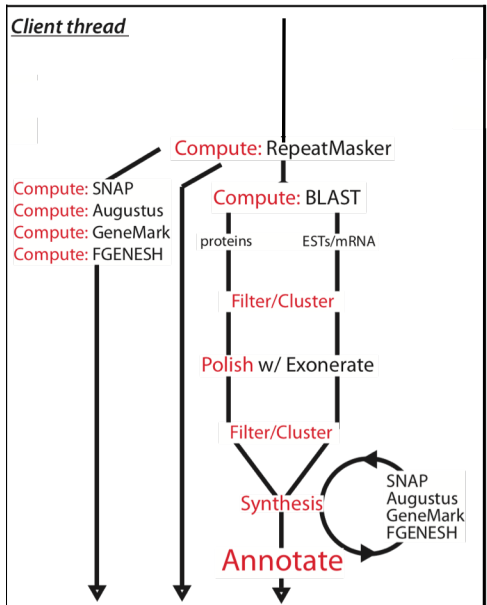
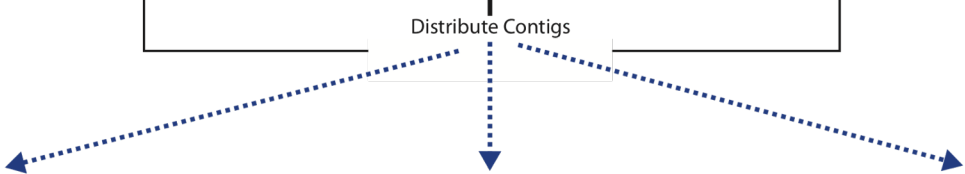


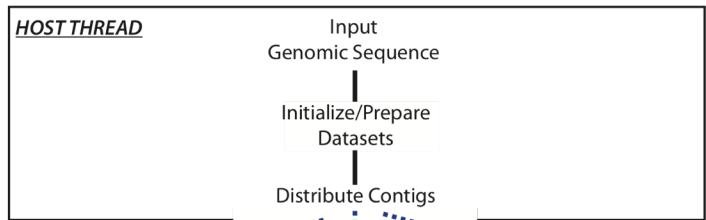




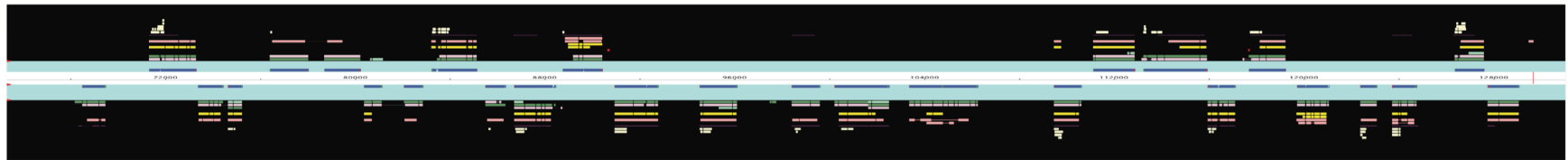
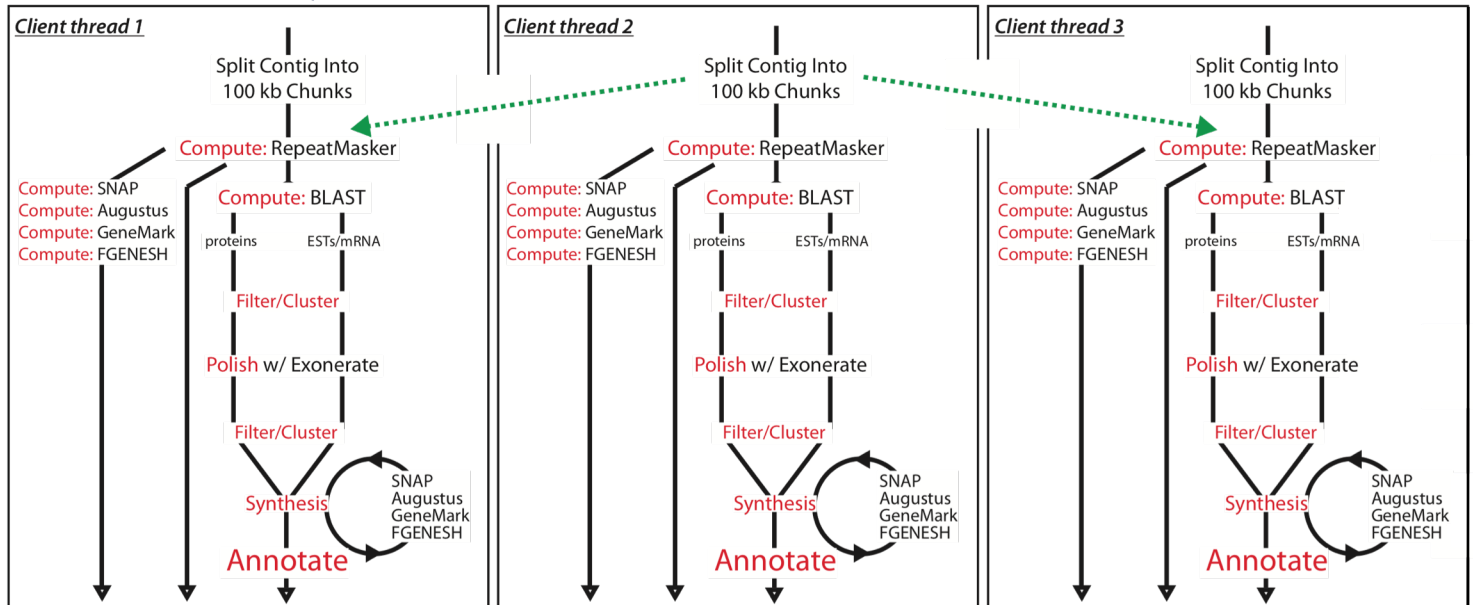


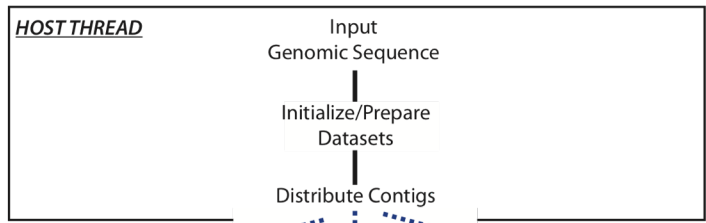
..... Contig level parallelization



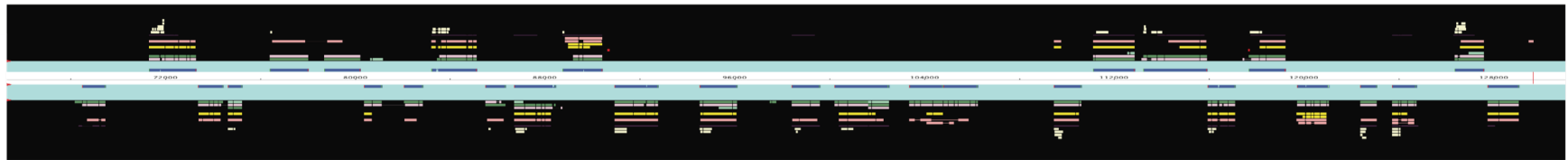
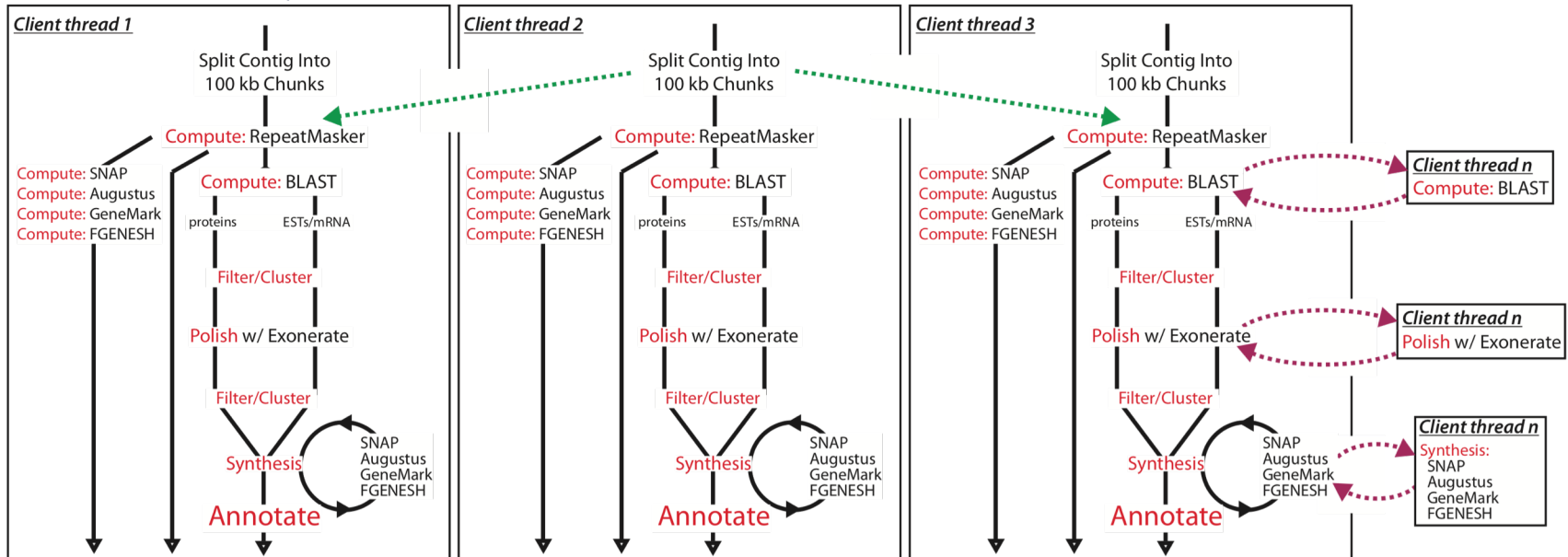


..... Contig level parallelization
 Parallelization by dividing sequence into chunks

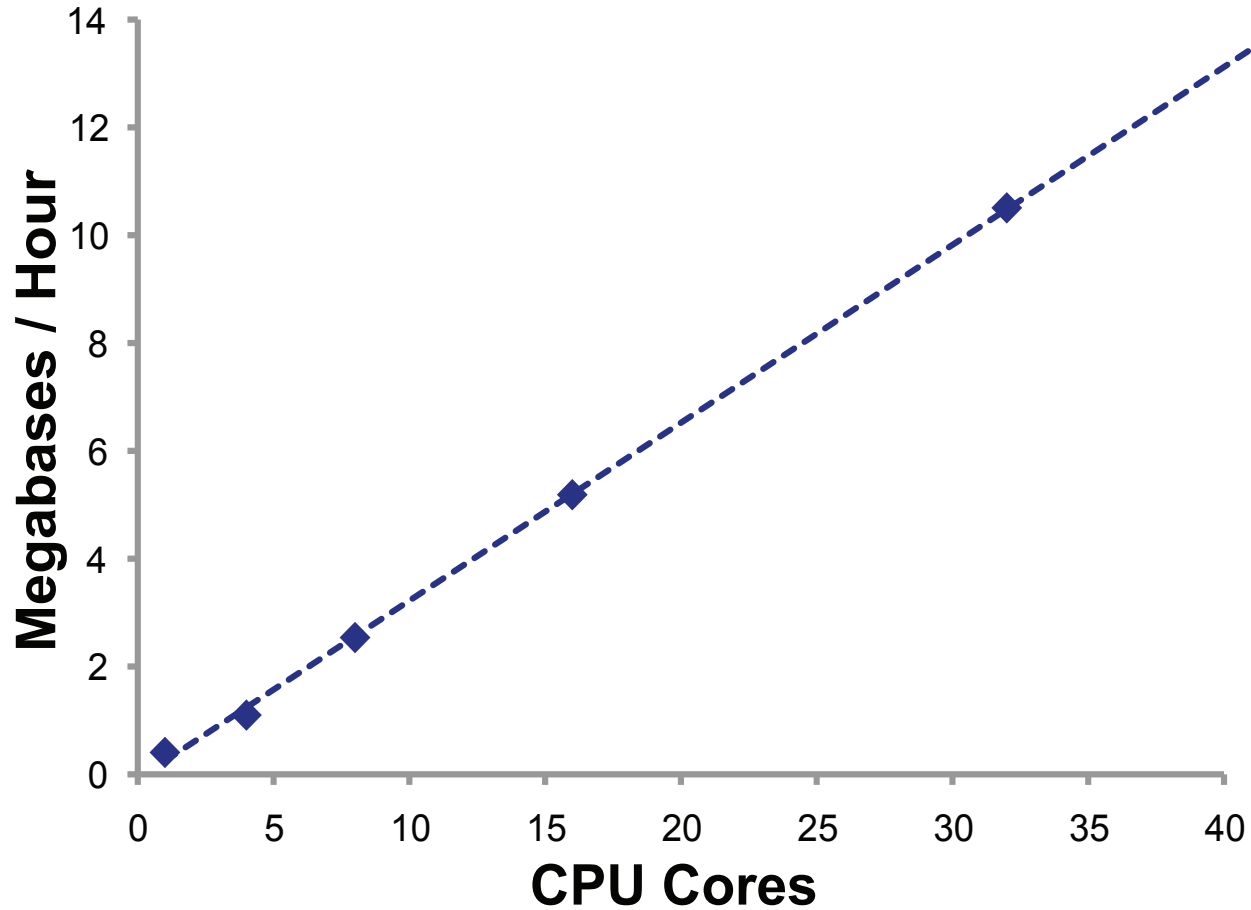




- Contig level parallelization
- Parallelization by dividing sequence into chunks
- Parallelization by distributing datasets among threads

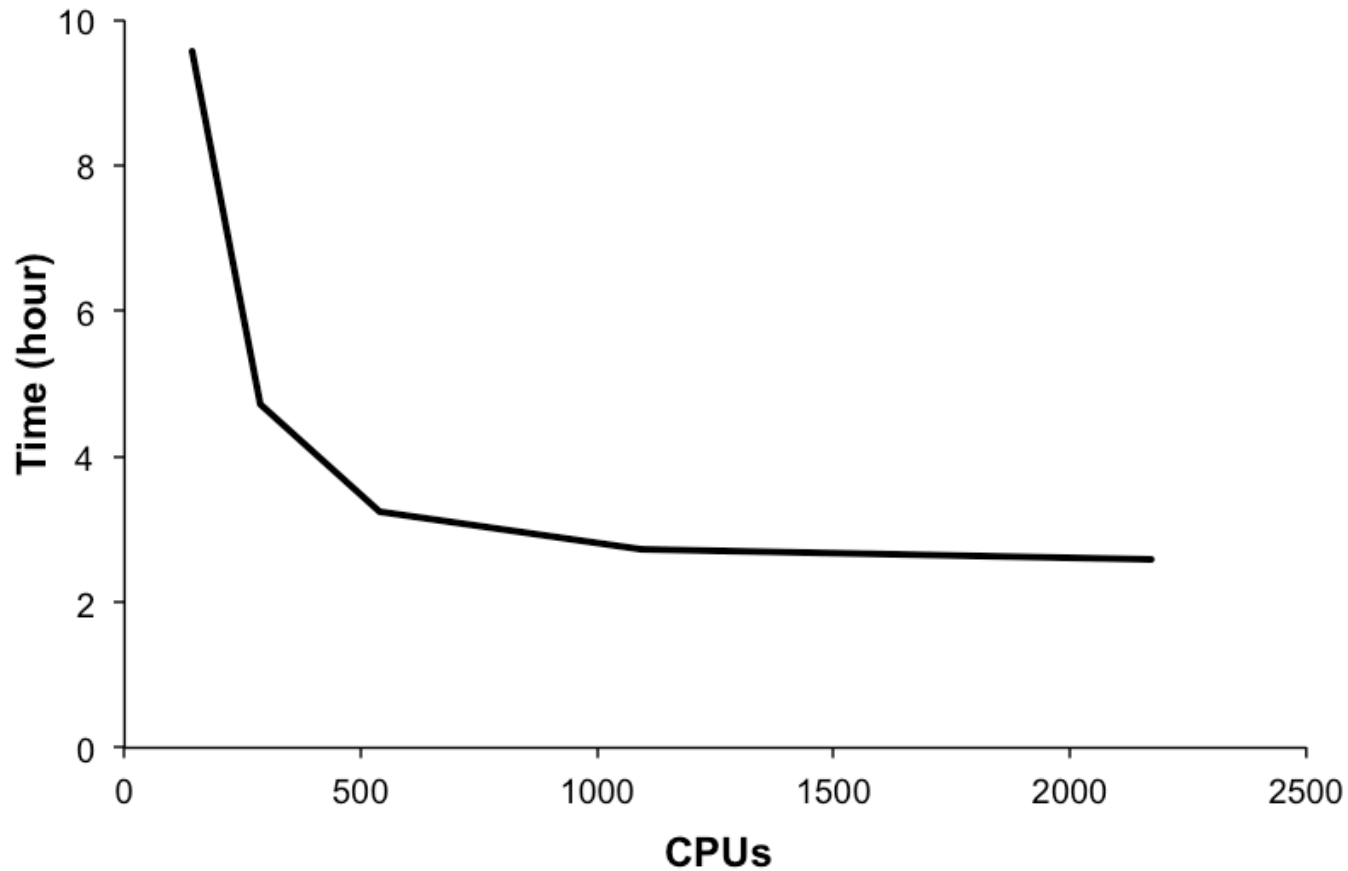


Data throughput

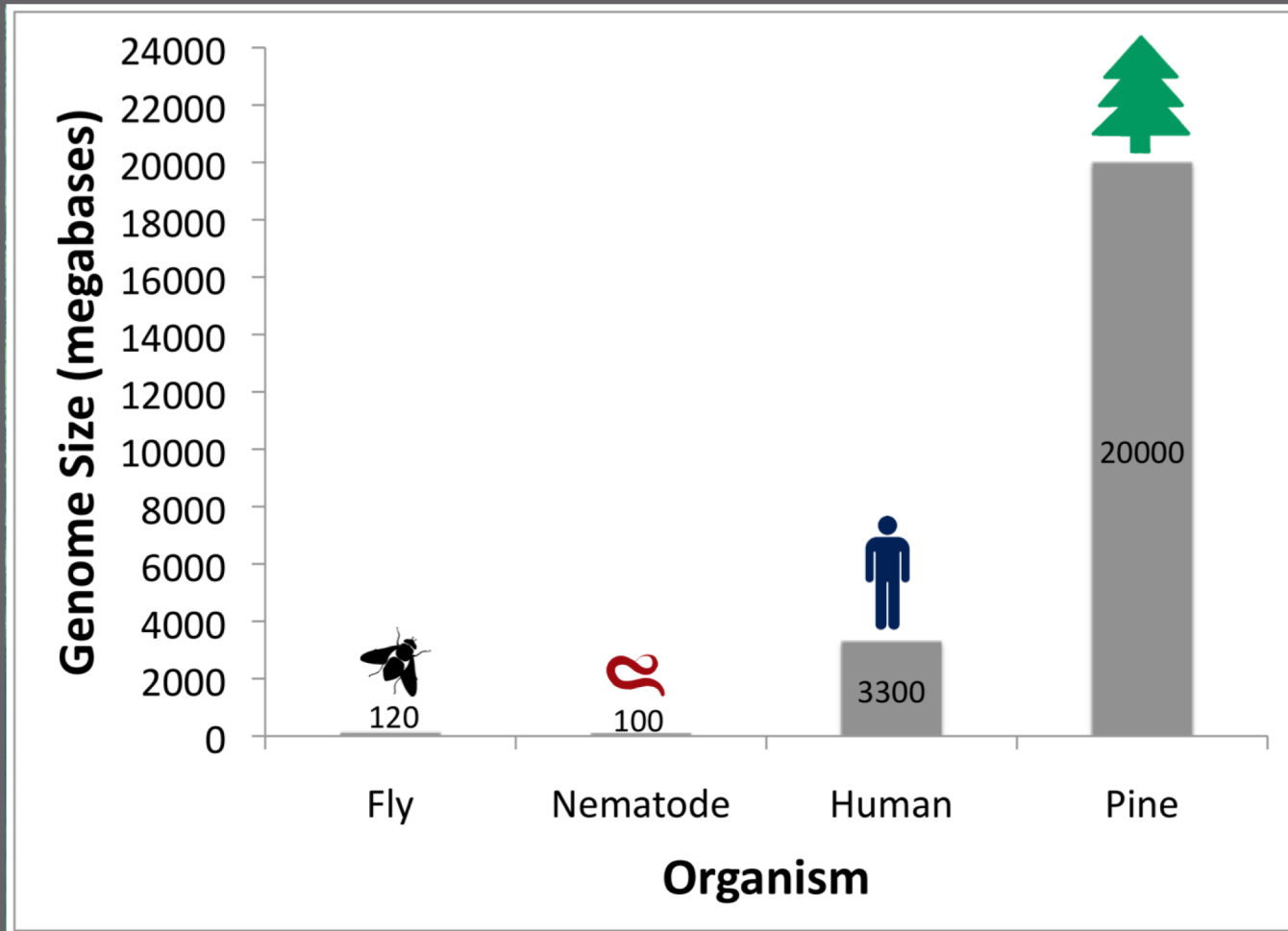


Data throughput

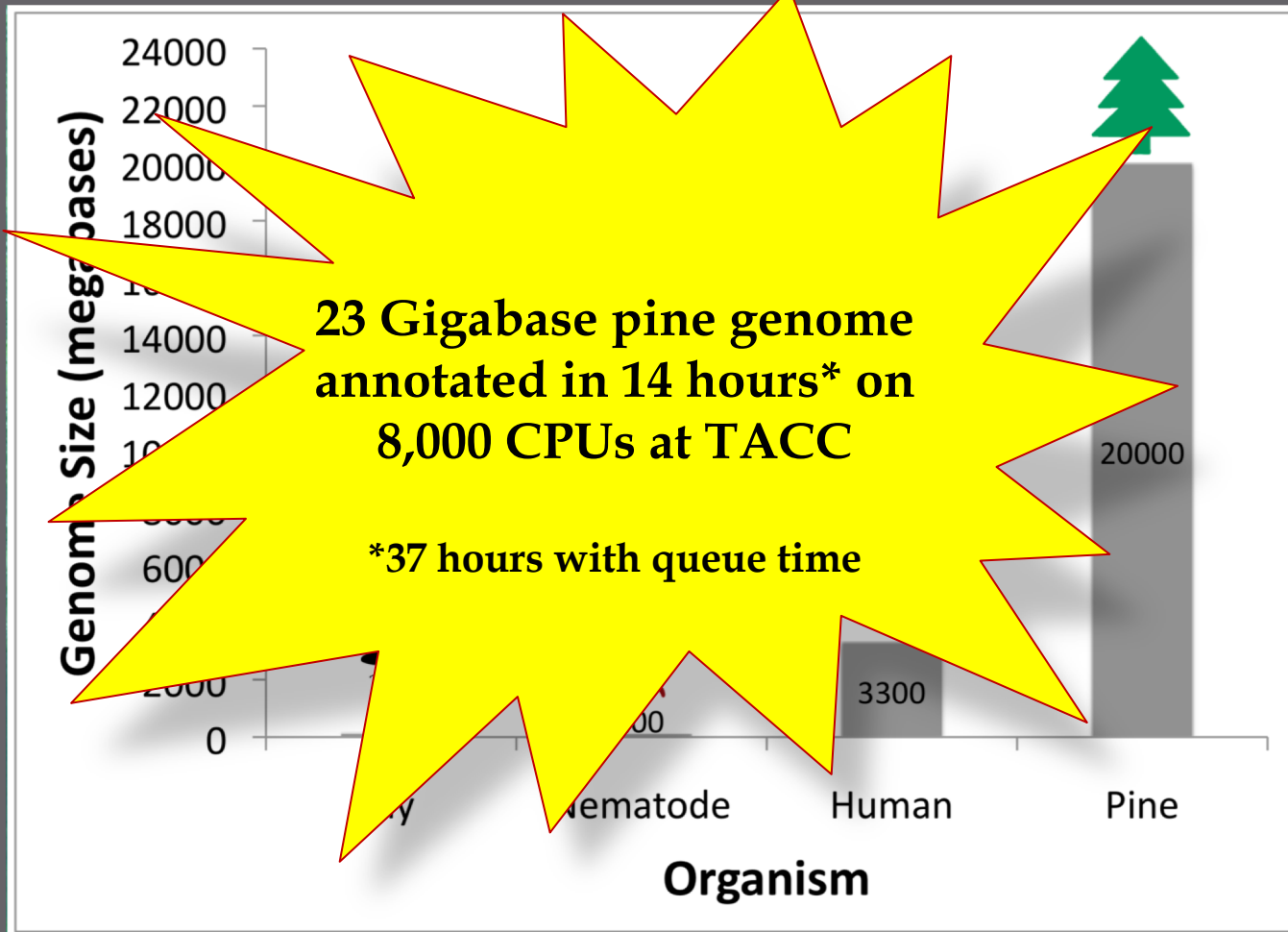
Annotation of the *Zea mays* Genome



Genome Sizes



Genome Sizes



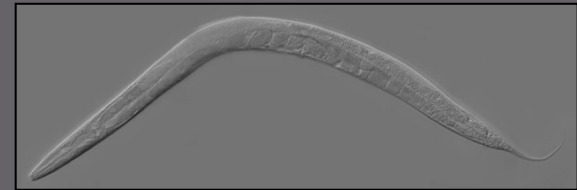
Next Generation of Genomics

“Second-generation ” genome projects have different needs and characteristics than earlier works.

First *versus* second-generation genomes

First-generation genomes:

- Classic experimental systems
- Large community
- Big \$
- Much prior knowledge about genome



Annotation Evaluation

Gene model accuracies for *ab initio* prediction and genome annotation programs

Reference Organism	Performance Category	<i>Ab Initio</i> Prediction			MAKER Annotations
		Augustus	GeneMark	SNAP	SNAP
<i>Caenorhabditis elegans</i>	Nucleotide Accuracy	88.29%	88.09%	85.10%	88.48%
<i>Arabidopsis thaliana</i>	Nucleotide Accuracy	77.04%	74.68%	69.78%	80.27%
<i>Drosophila melanogaster</i>	Nucleotide Accuracy	76.08%	66.54%	69.29%	74.33%

Annotation Evaluation

Gene model accuracies for *ab initio* prediction and genome annotation programs

Reference Organism	Performance Category	<i>Ab Initio</i> Prediction			MAKER Annotations
		Augustus	GeneMark	SNAP	SNAP
<i>Caenorhabditis elegans</i>	Nucleotide Accuracy	88.29%	88.09%	85.10%	88.48%
<i>Arabidopsis thaliana</i>	Nucleotide Accuracy	77.04%	74.68%	69.78%	80.27%
<i>Drosophila melanogaster</i>	Nucleotide Accuracy	76.08%	66.54%	69.29%	74.33%

With **enough** training data, *ab initio* gene predictors can match or even out-perform annotation pipelines*

*nGASP - the nematode genome annotation assessment project Avril Coghlan , Tristan J Fiedler , Sheldon J McKay ,Paul Flicek , Todd W Harris , Darin Blasiar , The nGASP Consortium and Lincoln D Stein BMC Bioinformatics 2008, 9:549doi:10.1186/1471-2105-9-549

First *versus* second-generation genomes

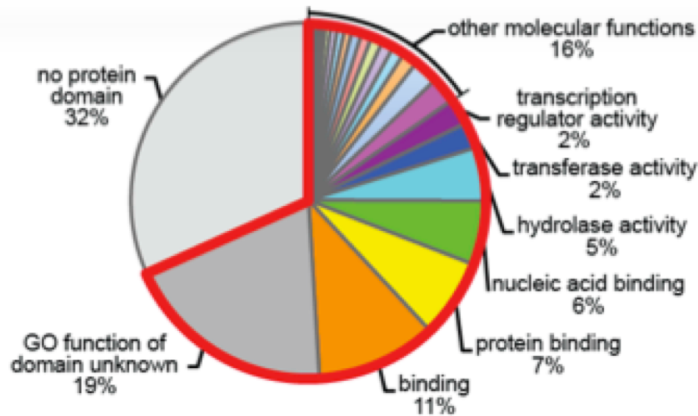
Second-generation genomes:

- New experimental systems
 - Genome will be the central resource for research
- Little prior knowledge about genome
 - Usually no genetics
- Small communities
- Less \$



Evaluation of *Ab initio* gene predictors on emerging model organism genomes

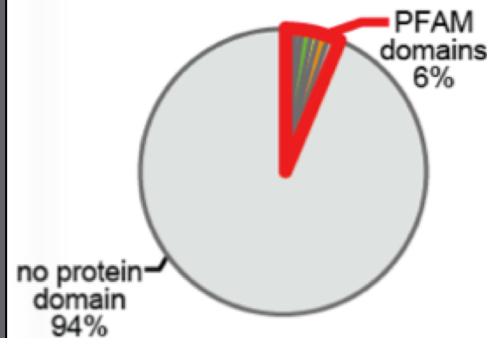
Average of Six Reference Proteomes



68% Contain Pfam Domain

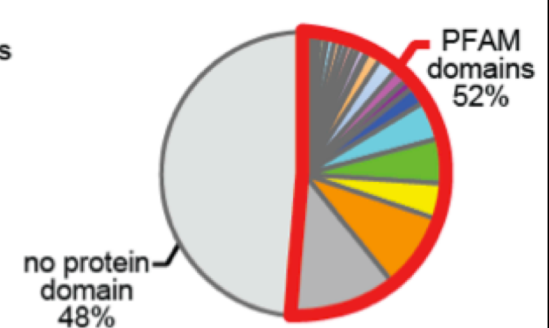
Schmidtea mediterranea

SNAP - *ab initio*



6% Contain Pfam Domain

MAKER - SNAP



52% Contain Pfam Domain

Evaluation of gene models with non-matched species parameters

Gene model accuracies for *ab initio* prediction and genome annotation programs

Reference Organism	Performance Category	<i>Ab Initio</i> Prediction		
		Augustus	GeneMark	SNAP
<i>Arabidopsis thaliana</i>	Nucleotide Accuracy	77.04%	74.68%	69.78%

Evaluation of gene models with non-matched species parameters

Gene model accuracies for *ab initio* prediction and genome annotation programs

Reference Organism	Performance Category	<i>Ab Initio</i> Prediction		
		Augustus	GeneMark	SNAP
<i>Arabidopsis thaliana</i>	Nucleotide Accuracy	77.04%	74.68%	69.78%

Gene model accuracies when using unmatched species parameter files

Reference Organism	Performance Category	<i>Ab Initio</i> Prediction		
		Augustus	GeneMark	SNAP
<i>Arabidopsis thaliana</i>	Nucleotide Accuracy	57.85%	48.62%	43.84%

Evaluation of gene models with non-matched species parameters

Gene model accuracies for *ab initio* prediction and genome annotation programs

Reference Organism	Performance Category	<i>Ab Initio</i> Prediction		
		Augustus	GeneMark	SNAP
<i>Arabidopsis thaliana</i>	Nucleotide Accuracy	77.04%	74.68%	69.78%
	Exon Accuracy	67.03%	61.31%	56.40%

Gene model accuracies when using unmatched species parameter files

Reference Organism	Performance Category	<i>Ab Initio</i> Prediction		
		Augustus	GeneMark	SNAP
<i>Arabidopsis thaliana</i>	Nucleotide Accuracy	57.85%	48.62%	43.84%
	Exon Accuracy	30.71%	16.51%	18.58%

Evaluation of gene models with non-matched species parameters

Gene model accuracies for *ab initio* prediction and genome annotation programs

Reference Organism	Performance Category	<i>Ab Initio</i> Prediction		
		Augustus	GeneMark	SNAP
<i>Arabidopsis thaliana</i>	Nucleotide Accuracy	77.04%	74.68%	69.78%
	Exon Accuracy	67.03%	61.31%	56.40%

Gene model accuracies when using unmatched species parameter files

Reference Organism	Performance Category	<i>Ab Initio</i> Prediction		
		Augustus	GeneMark	SNAP
<i>Arabidopsis thaliana</i>	Nucleotide Accuracy	57.85%	48.62%	43.84%
	Exon Accuracy	30.71%	16.51%	18.58%

Evaluation of gene models with non-matched species parameters

Gene model accuracies for *ab initio* prediction and genome annotation programs

Reference Organism	Performance Category	<i>Ab Initio</i> Prediction		
		Augustus	GeneMark	SNAP
<i>Arabidopsis thaliana</i>	Nucleotide Accuracy	77.04%	74.68%	69.78%
	Exon Accuracy	67.03%	61.31%	56.40%

Gene model accuracies when using unmatched species parameter files

Reference Organism	Performance Category	<i>Ab Initio</i> Prediction			MAKER Annotations
		Augustus	GeneMark	SNAP	SNAP
<i>Arabidopsis thaliana</i>	Nucleotide Accuracy	57.85%	48.62%	43.84%	73.77%
	Exon Accuracy	30.71%	16.51%	18.58%	60.11%

Evaluation of gene models with non-matched species parameters

Gene model accuracies for *ab initio* prediction and genome annotation programs

Reference Organism	Performance Category	<i>Ab Initio</i> Prediction		
		Augustus	GeneMark	SNAP
<i>Arabidopsis thaliana</i>	Nucleotide Accuracy	77.04%	74.68%	69.78%
	Exon Accuracy	67.03%	61.31%	56.40%
<i>Drosophila Melanogaster</i>	Nucleotide Accuracy	76.08%	66.54%	69.29%
	Exon Accuracy	61.37%	47.31%	47.01%

Gene model accuracies when using unmatched species parameter files

Reference Organism	Performance Category	<i>Ab Initio</i> Prediction			MAKER Annotations
		Augustus	GeneMark	SNAP	SNAP
<i>Arabidopsis thaliana</i>	Nucleotide Accuracy	57.85%	48.62%	43.84%	73.77%
	Exon Accuracy	30.71%	16.51%	18.58%	60.11%
<i>Drosophila Melanogaster</i>	Nucleotide Accuracy	67.47%	66.51%	48.92%	74.44%
	Exon Accuracy	30.62%	26.25%	19.94%	53.69%

Evaluation of gene models with non-matched species parameters

Gene model accuracies for *ab initio* prediction and genome annotation programs

Reference Organism	Performance Category	<i>Ab Initio</i> Prediction		
		Augustus	GeneMark	SNAP
<i>Arabidopsis thaliana</i>	Nucleotide Accuracy	77.04%	74.68%	69.78%
	Exon Accuracy	67.03%	61.31%	56.40%
<i>Drosophila Melanogaster</i>	Nucleotide Accuracy	76.08%	66.54%	69.29%
	Exon Accuracy	61.37%	47.31%	47.01%
<i>Caenorhabditis elegans</i>	Nucleotide Accuracy	88.29%	88.09%	85.10%
	Exon Accuracy	74.62%	68.88%	61.38%

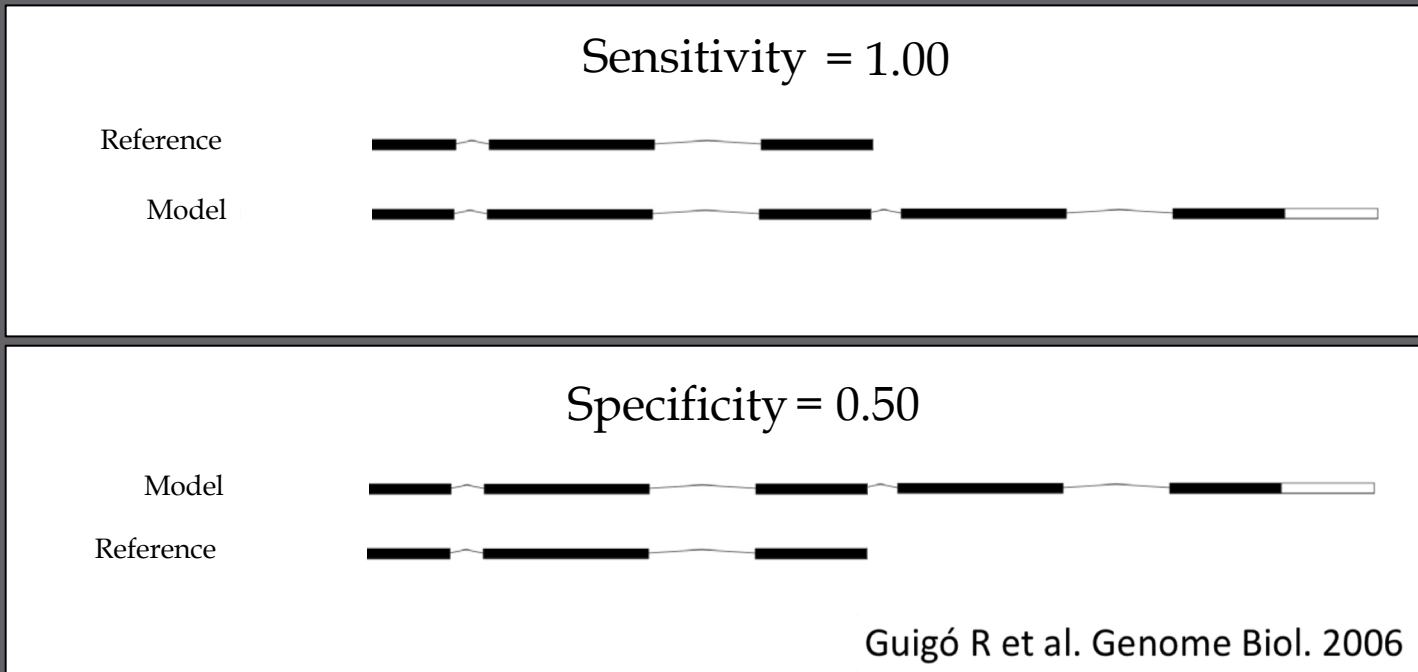
Gene model accuracies when using unmatched species parameter files

Reference Organism	Performance Category	<i>Ab Initio</i> Prediction			MAKER Annotations
		Augustus	GeneMark	SNAP	SNAP
<i>Arabidopsis thaliana</i>	Nucleotide Accuracy	57.85%	48.62%	43.84%	73.77%
	Exon Accuracy	30.71%	16.51%	18.58%	60.11%
<i>Drosophila Melanogaster</i>	Nucleotide Accuracy	67.47%	66.51%	48.92%	74.44%
	Exon Accuracy	30.62%	26.25%	19.94%	53.69%
<i>Caenorhabditis elegans</i>	Nucleotide Accuracy	66.18%	67.26%	68.24%	85.02%
	Exon Accuracy	28.33%	30.01%	35.44%	63.14%

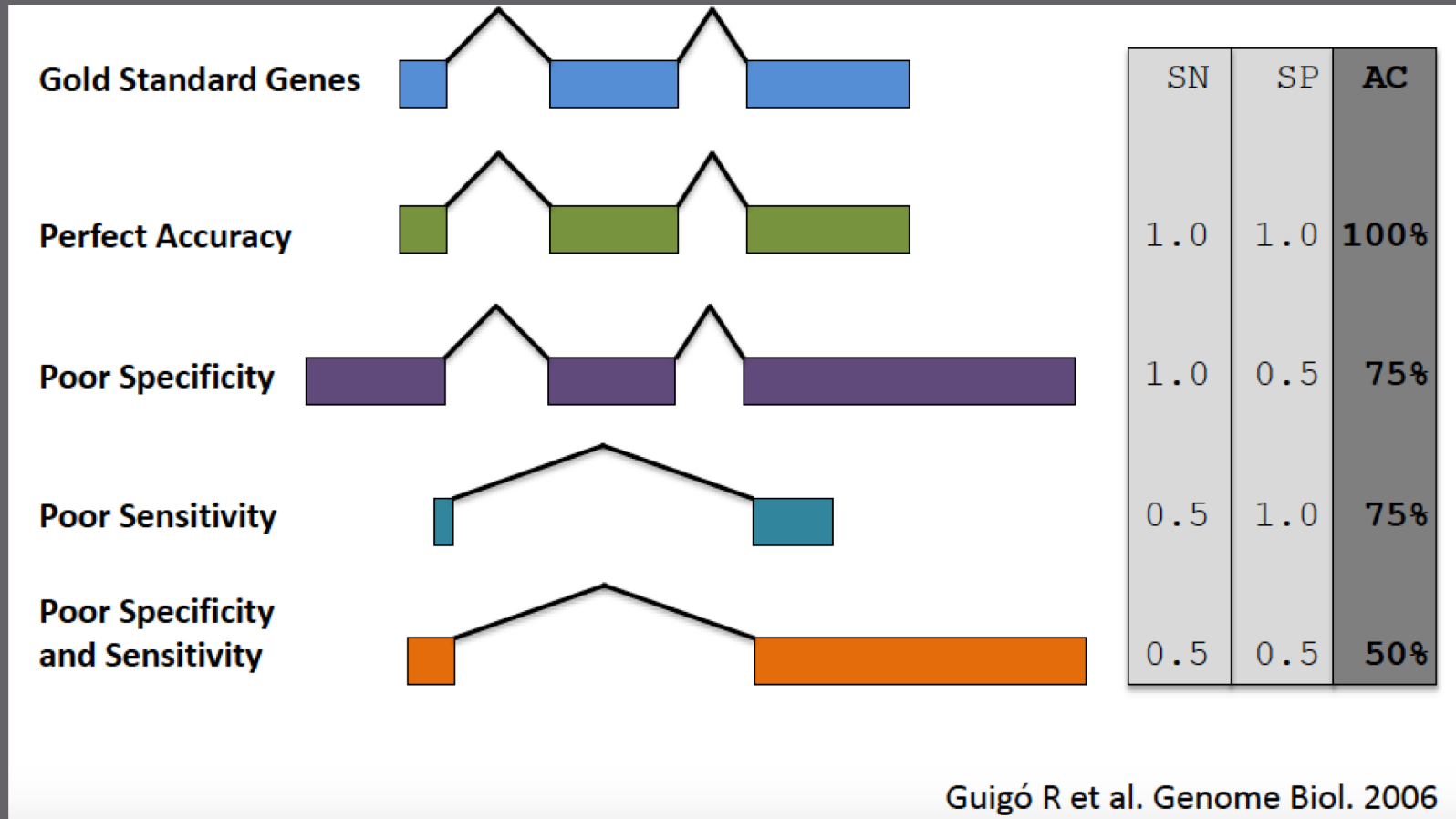
Beyond *de novo* annotation

- Quality control and data prioritization
- Update/revise legacy annotation sets
- Integrating new evidence into existing databases

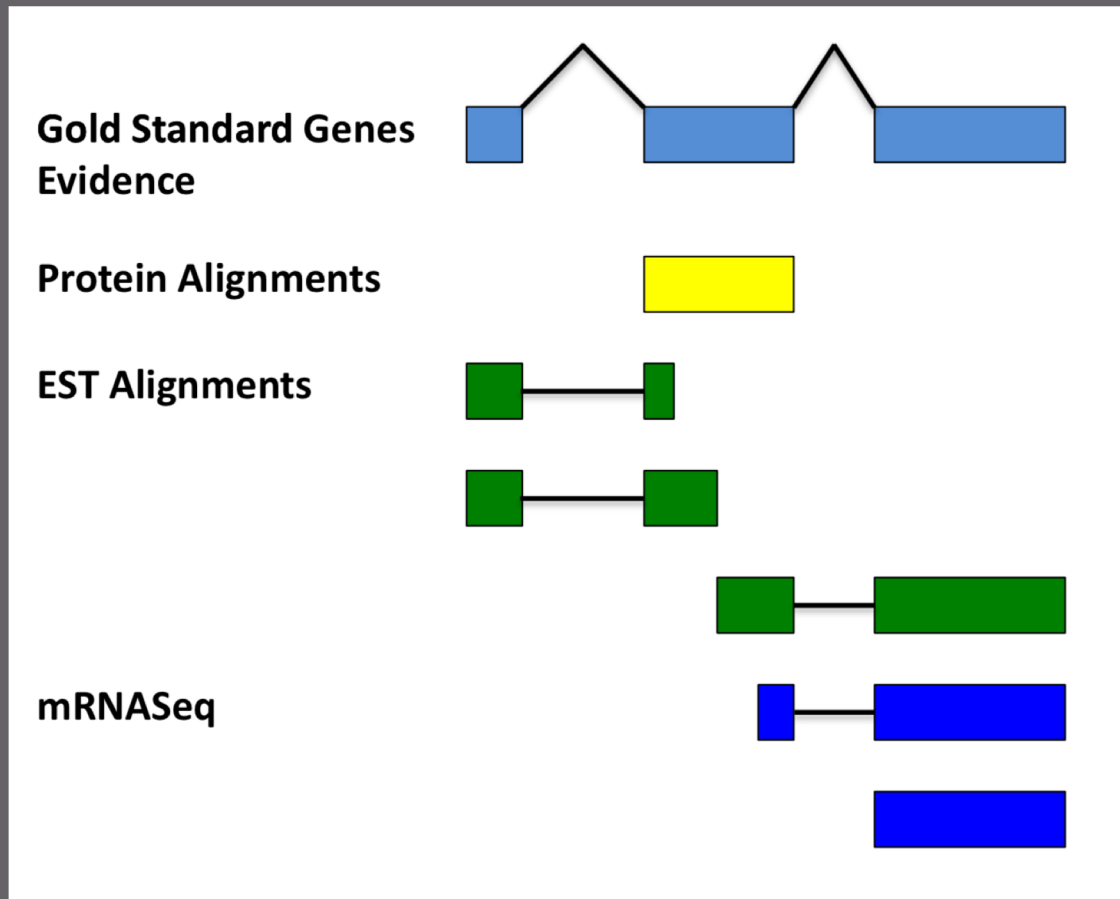
Quality control and data prioritization



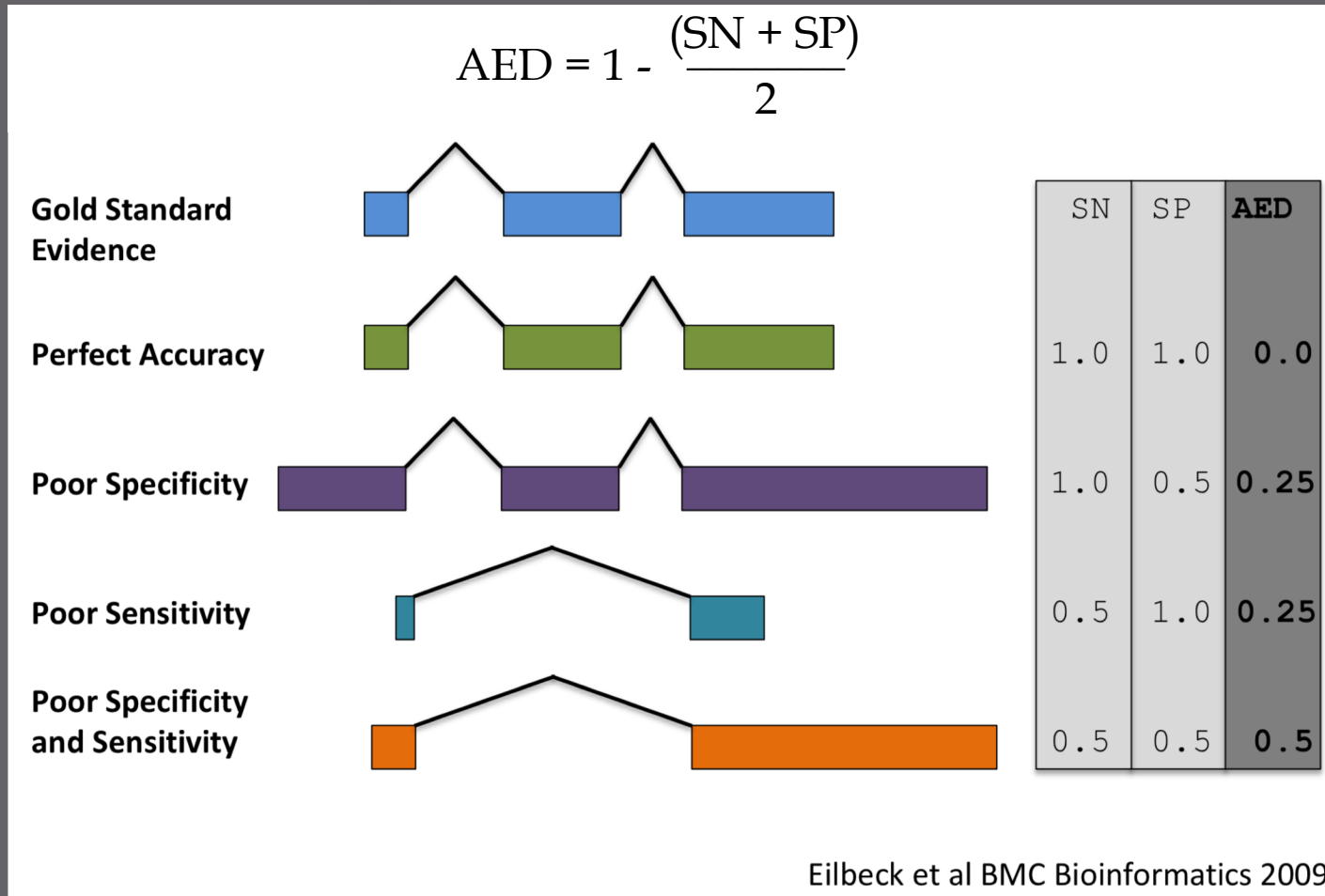
Quality control and data prioritization



Quality control and data prioritization

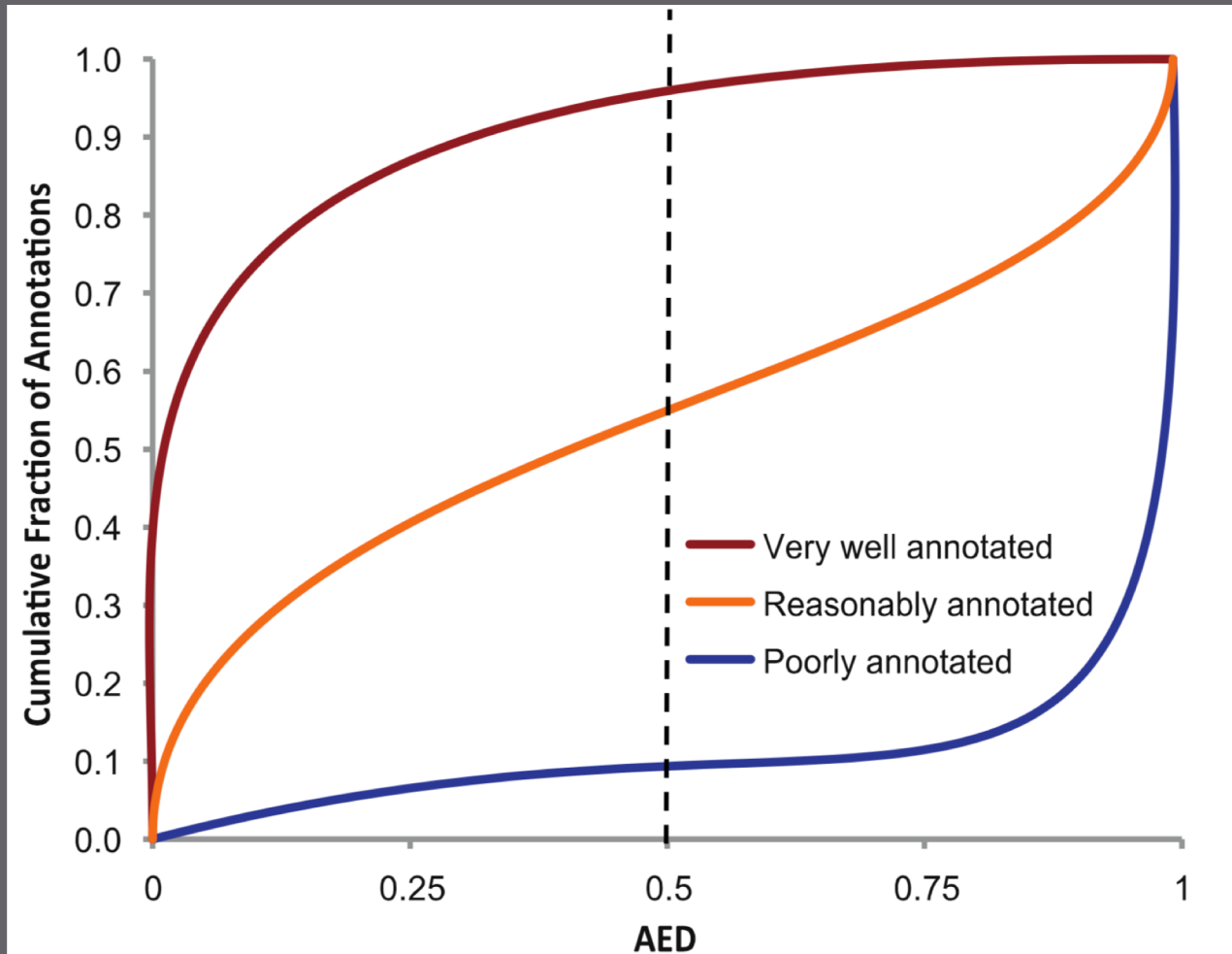


Quality control and data prioritization



*Quantitative Measures for the Management and Comparison of Annotated Genomes
 Karen Eilbeck , Barry Moore , Carson Holt and Mark Yandell BMC Bioinformatics 2009
 10:67doi:10.1186/1471-2105-10-67

Quality control and data prioritization



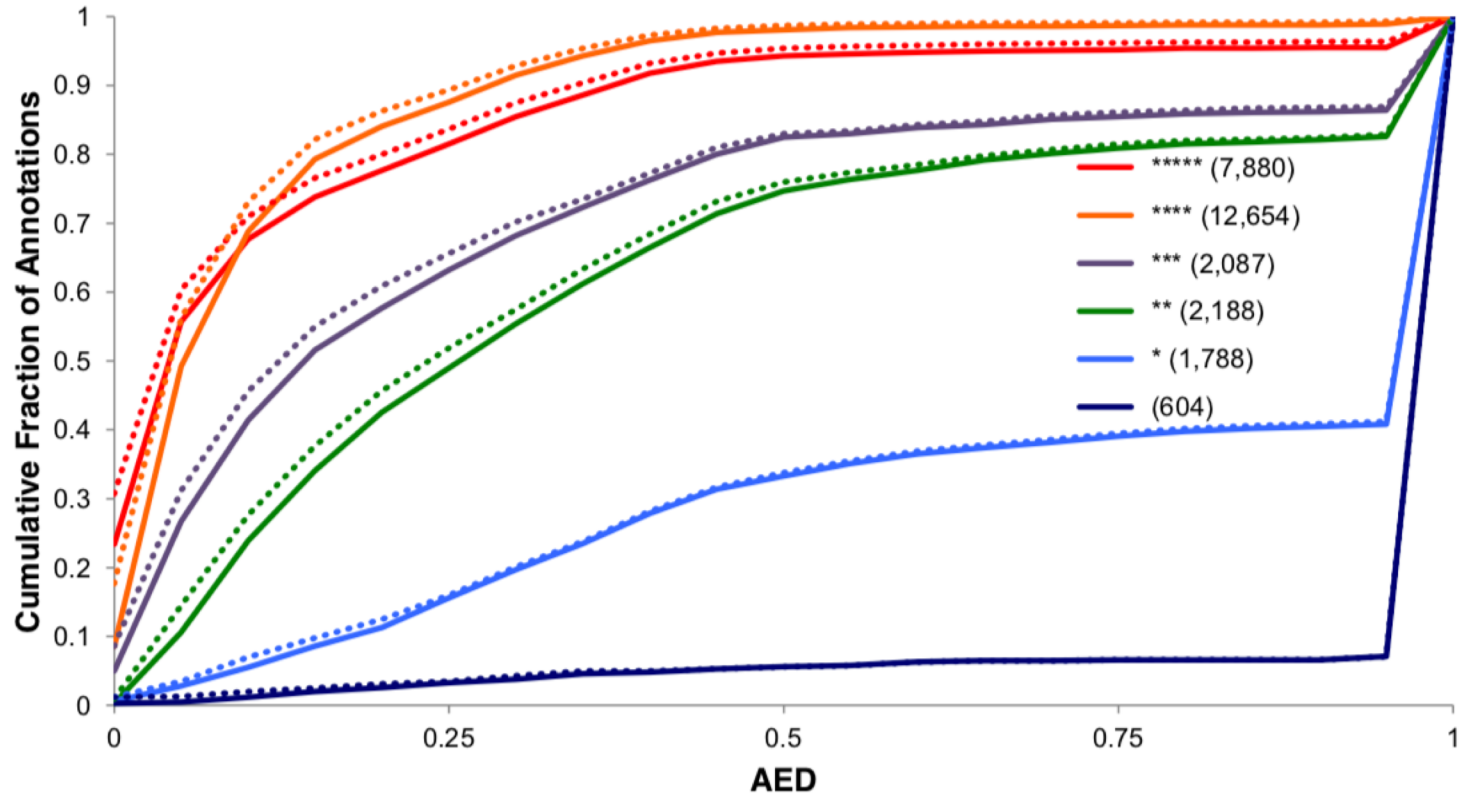
Evaluation of AED on Arabidopsis

TAIR Star Rating System

Gene Models	Internal Exons	External Exons	Single Exon
1 ***** Every Splice junction * & complete coverage ** by a single piece of evidence	E1 & all ex's in same transcript	X1	S1
2 ***** Every splice junction all aligning evidence consistent Complete Coverage**	E1 E2	X1	
3 ***** Every splice junction some aligning evidence inconsistent Complete Coverage**	E1 E2	X1	
4 ***** Every splice junction all aligning evidence consistent Not Complete Coverage	E1 E2 & E3	X1 X2 X3	S2
5 ***** Every Splice junction some aligning evidence inconsistent Not complete Coverage	E1 E2 & E3	X1 X2 X3	
6 ***** Not every Splice junction >50% coverage	E1 E2 E3 & E4 E5	any X	
7 ***** Not every Splice junction	E1 E2 E3 & E4 E5	any X	
8 ***** No Splice junctions >50% coverage	E5	X4	S3
9 ***** No Splice junctions	E5	X4	S3
10 ***** No experimental evidence	E6	X5	S4

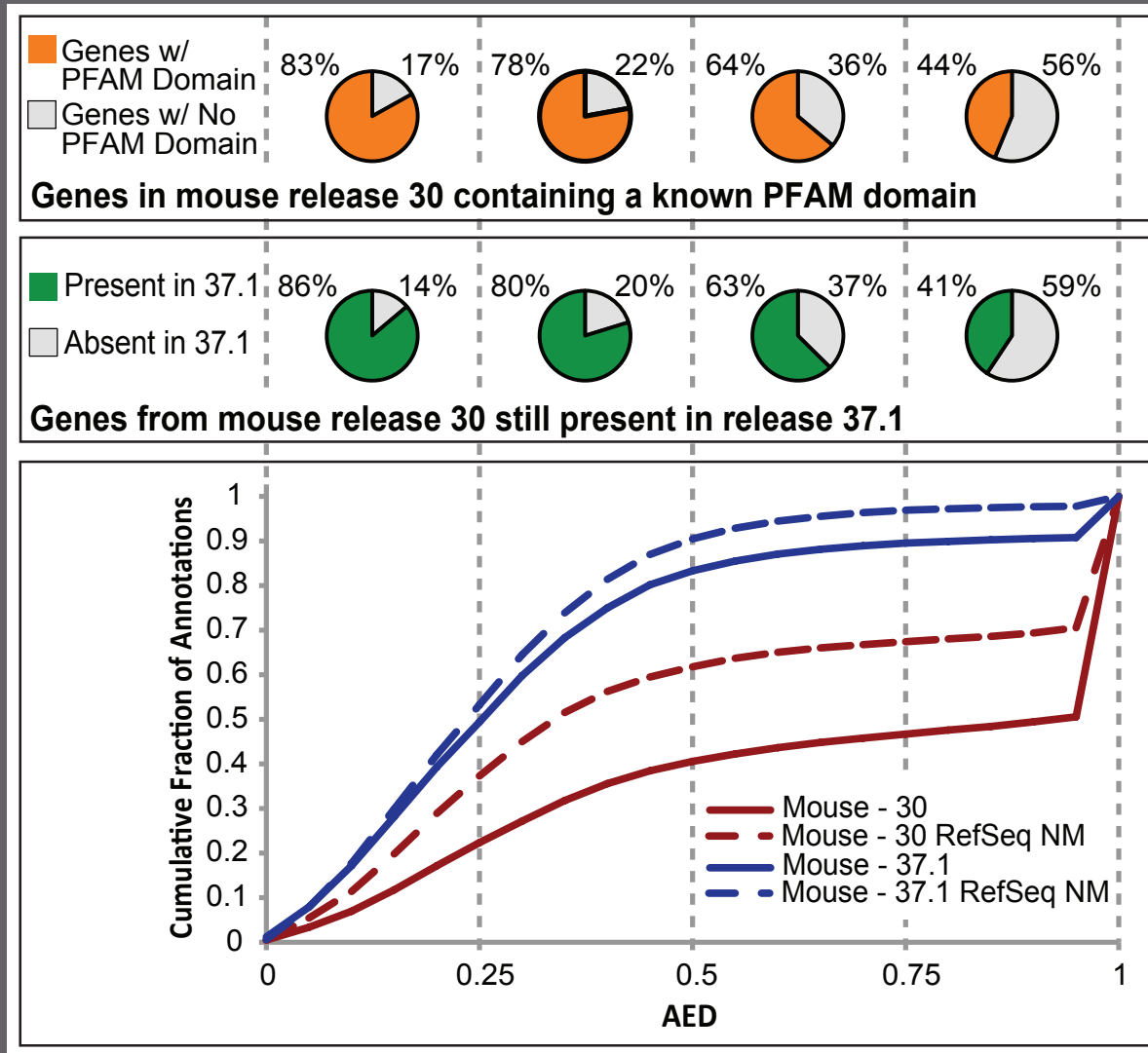
*except for single exon genes
 ** Transcript may extend beyond end of genemodel

Evaluation of AED on Arabidopsis

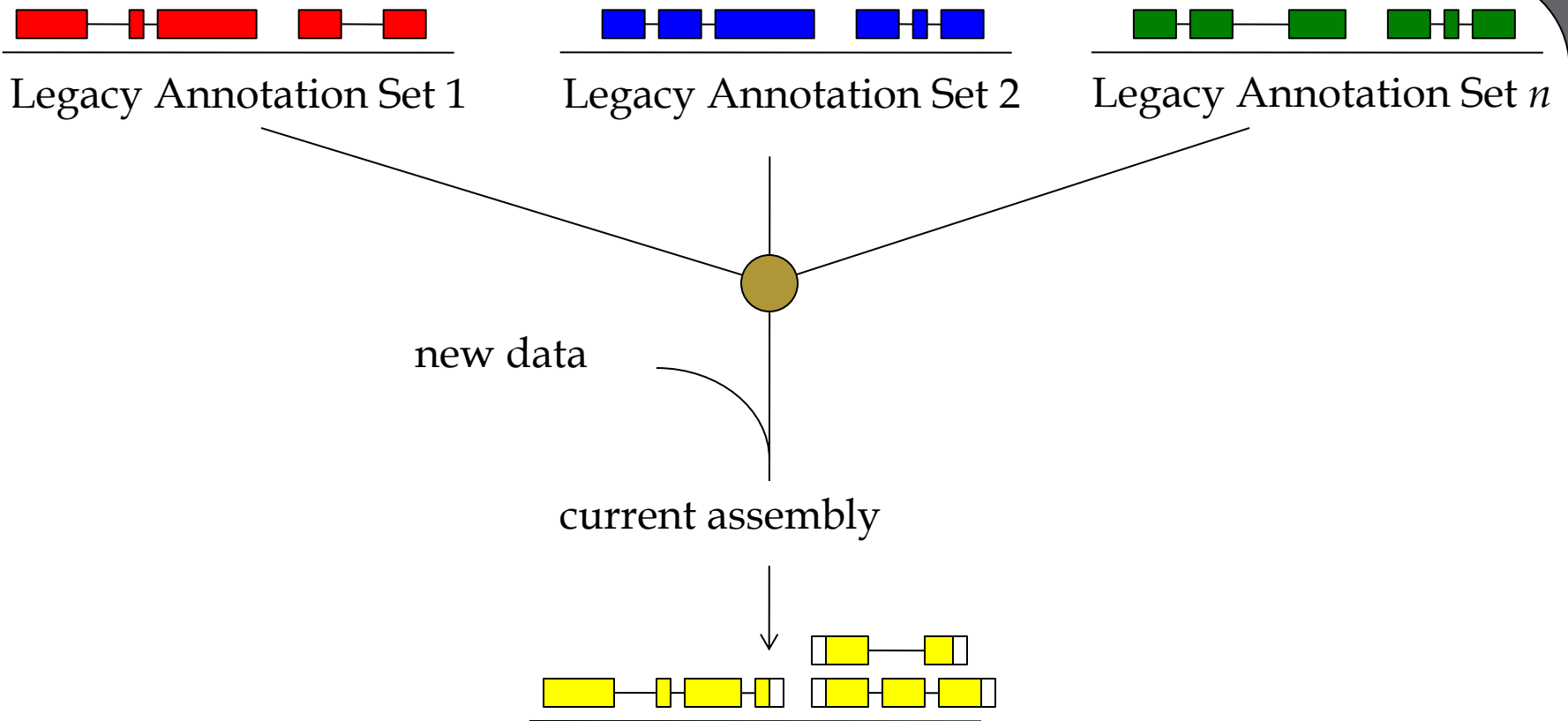


Evidence: mRNA-seq (17 experiments), ESTs, full length cDNAs, Swiss-Prot (minus Arabidopsis)

Evaluation of AED on the mouse genome

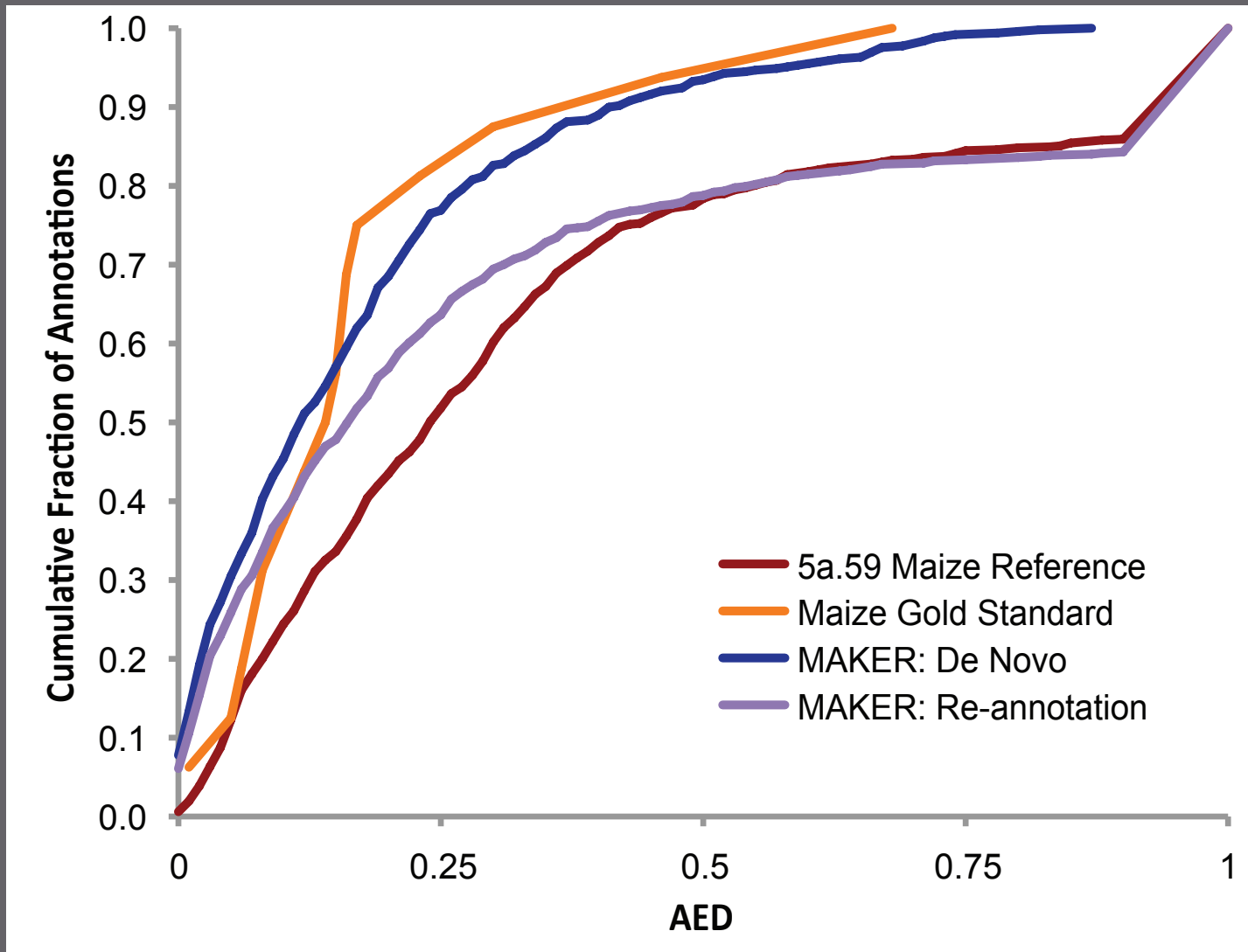


Update/revise legacy annotation sets



- Identify legacy annotation most consistent with new data
- Automatically revise it in light of new data
- If no existing annotation, create new one

Re-annotation of the maize genome



MAKER Versions

▣ MAKER

- 2008. Based on earlier annotation pipelines developed by Mark Yandell

▣ MAKER 2 / MAKER-P

- 2011. Introduction of MPI parallelization, support for multiple gene predictors, GFF3 pass-through, and quality metrics like AED (Annotation Edit Distance) from the Sequence Ontology consortium.
- 2015. Support for tRNA and snoRNA annotation. Improved parallelization on large plant genomes.

▣ MAKER 3

- 2016. EVM (Evidence Modeler) support for improved annotation and user defined evidence probability weighting.



MAKER Web Annotation Service

Your Genome Annotated

- Home
- Help
- Yandell Lab

not logged-in | [sign in](#)

Welcome to the MAKER Web Annotation Service

Log into your account below, or you can access the server as a guest. While there is no login requirement for this site, users are highly encouraged to set up an account. Use the "New user registration" link to register a new account. Registration is free, and has several benefits. Registered users can submit up to 5,000,000 base pairs of sequence for each annotation job. Guest users are limited to 500,000 base pairs per annotation job submission. Registered users

<http://www.yandell-lab.org/>



Maker Web Annotation Service

User Name

Password

Remember User Name

- [New user registration](#)
- [Forgot login?](#)
- [Help](#)

New Guest Account

User Sign In



Service, click "Help" above.

Example jobs are meant to be used with the [step by step tutorial](#) hosted by the Generic Model Organism Database Project. After clicking load for an example job, the fields below will be filled out for you. You can then review them or edit them. You will need to click on "Add Job to Queue" at the bottom of the page before the example job will start.

Example Jobs

 MAKER Job Details Assigned id: 1656

```
>NT_010783.15 section 6410321-6611764 Homo sapiens chromosome 17 genomic contig, GRCh37 reference primary asse
TTTTCTGTAGAGGGTGGGGCTCTCCCTATGTTGCTCAGCTGGTCTCTGACTCCTGGGCTCAAGCCACCCCTCC
CACCTTAGCCTCCCTAAGTGTGGGATTACAGGCATGAGCCACTGCACCTGGCCCCAGTTTTTTTTTTT
TTTTTTTGAGACAGAGTCTGGCCCTGCACCCAGGCTGGAGTGCACTAGCATGATCTCAGCTCACGCAA
CCTCCACCTCCTAGTTCAAGCTGTTCTGCTGCCTCAGCTCCCTAGTAGCTGGGACTGCAGGCATGTGC
CACTATCCCCGGCTAAAATTTGTATTTTGTAGTAGAGACAGGGTTTGGCCATGTTGGCCAGGCTGGTCTCG
AACTCCTGACCTCAAGTGATCCACCCACCTCAACCCCAAATATTTTAAACCAATATTTCTCTTATCTTA
CTTCATACTGGAATGATTTATGTTATTTCCTGCTTGCACCAGAGCTCATTGTTCTCTTGTTTGGCC
ACACTCCAGCAAAATAGCATATCGTAGCTACCATGGCAAGTGACAGAGTCAACAGACTTGGTCATGAAAA
TGGAAATCAAAATTTATGGCTGATTTCTTCTCCCTGACATAAAGTTGGGTAACAGTGAGAACAAAAACA
AAAGACTTTTCCACAGTTCTTCTCCTTAGGCTGGTTAGACAACTTCTCCACTCAGAGTTCATTATTTG
GAAATTTGTCATTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTT
GTCCTGCTGTGACACCAGGCTGGAGTGCAATGCTGCGGTCTTGGCTCAGTCCAACTCCGCTCCTGGG
TTCAGCCATTTCTCCTGCTCAGCTCCCAAGTAGCTAGGATTACAGGCACGTGCCACCACACTGGGCTA
GTTTTGTAGAGACGGGGTTCCACCATGTTGGCCAGGCTGGTCTCGAACTCCTGACCTCAGGTGATCTGC
CTGCTTTGGCCTCCCAAAGTGTGGGATTACAGGCATGAGCCACTGCGTTGGCCATACTTCTCTAATT
CAGTTCAGTAACATACAACAAATGCATGTTGAGAATTAATCTTGTGCCAGACACTGTGCTGGCATTGTT
GGTTCAAAGTCAAAACCTTTGCTCTTAAGGAGTTCATAGACTAGAAGAGACAGAAATCCATAACAAATAAT
GTGATATGAGCTAGACTAGAAGTGTGAACAACCTCTAGAGGAAATGTTTAGTTATGCTGGGAAATCAG
GAAGCCTCTTAAGGAGAAGTAGGCTTCCACCAATGGGGAAGTGGGATAGGTTCTAACAGTATGCAAC
TGCCTCTGAAAGGATTACAGAGGACTGCCCTGAGAAAAGGCTTGCCTTCATGTTGAGCAAAGTTAAAG
CAAGAGGATTGCATTTTAAATACTTAAACAAAAAGTTAATTTGTTGTTTTCTTTTTTTTTTTTTTTTCAT
TTCCTCTCTCATTACTTACTTCTCAGGGCTTGGTAAGTACCGTTTAAAGTTAAATGTAATACGCTG
TAGACTTCTAGTTTAAAGACGAGCCTTGGCCCTACAATGATTTGAGGTTTTGCTTTTTGTTTTGAGA
CTCAAGAAAGACAAGGAAGAACCAGTGTGAGAATACGTTGACAAACTTTACCACCTGAAGCATTTT
TATTAGTTGAAATGCAGAACTACTGGTGGCCAAACAGTTGAAAAAGGAGTCAGTCTCAGCAAACTTAC
AGATACTTTTTTTTAAAAACGTTGTTACCTTACATCCCAATTATGTTGGCTTTTGGTTTTTTTTTCTCTGG
```

ESTs from the same source as your genomic sequence; and (2) (optionally) ESTs from a closely related organism, for example if your genomic sequence is human, this second set of ESTs might be from mouse.

Upload a multi-fasta file of ESTs to be aligned from the same source

No file selected...



MAKER Web Annotation Service

Your Genome Annotated

- Home
- New Job
- Manage Files
- Running Jobs
- Edit Account
- Contact Us
- Help
- Yandell Lab

logged-in as [guest_270](#) | [logout](#)

You are logged in as a guest user. You may become a registered user at any time by clicking on "Edit Account".

Copy the URL <http://derringer.genetics.utah.edu/cgi-bin/MWAS/maker.cgi> or login with the username [guest_270](#) to return to the main page.

Annotation Status Summary:

TOTAL Contigs: 1
 FINISHED: 1
 INCOMPLETE: 0
 FAILED: 0
 SKIPPED: 0

[Download All Data](#)
 Do post processing of annotations

View contigs individually.

- [View in GBrowse](#)
- [View in JBrowse](#)
- [View in Apollo](#)
- [SOBA Statistics](#)

Welcome to the MAKER

To get started just click on "New Job" above. You can also click on "Example Annotations" to see pre-loaded example annotations as the results below. You can also click on "Job Queue" above. For more information see the "Help" page.

[Refresh Job Status](#)

Your Jobs (1)

JobID	Description
1323	Chromosome 17

Clicking on "Launch in Apollo" will install a Java Web Start version of Apollo if not already installed. If for some inexplicable reason the program Dashboard starts [click here](#).



MAKER Web Annotation Service: 201.4 kbp from NT_010783.15:1..201,444
http://derringer.genetics.utah.edu/cgi-bin/gb2/gbrowse/MWAS_270_1323/?name=NT_010783.15

File Help

Browser Select Tracks Upload and Share Tracks Preferences

Search
Landmark or Region: NT_010783.15:1..201,444 Search
Annotate Restriction Sites [v] Configure... Go
Data Source: MAKER Web Annotation Service [v]
Scroll/Zoom: [<<] [<] [] [>] [>>] Show 201.4 kbp [v] Flip []

Overview
NT_010783.15
0k 10k 20k 30k 40k 50k 60k 70k 80k 90k 100k 110k 120k 130k 140k 150k 160k 170k 180k 190k 200k

Region
0k 10k 20k 30k 40k 50k 60k 70k 80k 90k 100k 110k 120k 130k 140k 150k 160k 170k 180k 190k 200k

Details
NT_010783.15: 201.4 kbp
50 kbp |-----|
0k 10k 20k 30k 40k 50k 60k 70k 80k 90k 100k 110k 120k 130k 140k 150k 160k 170k 180k 190k 200k

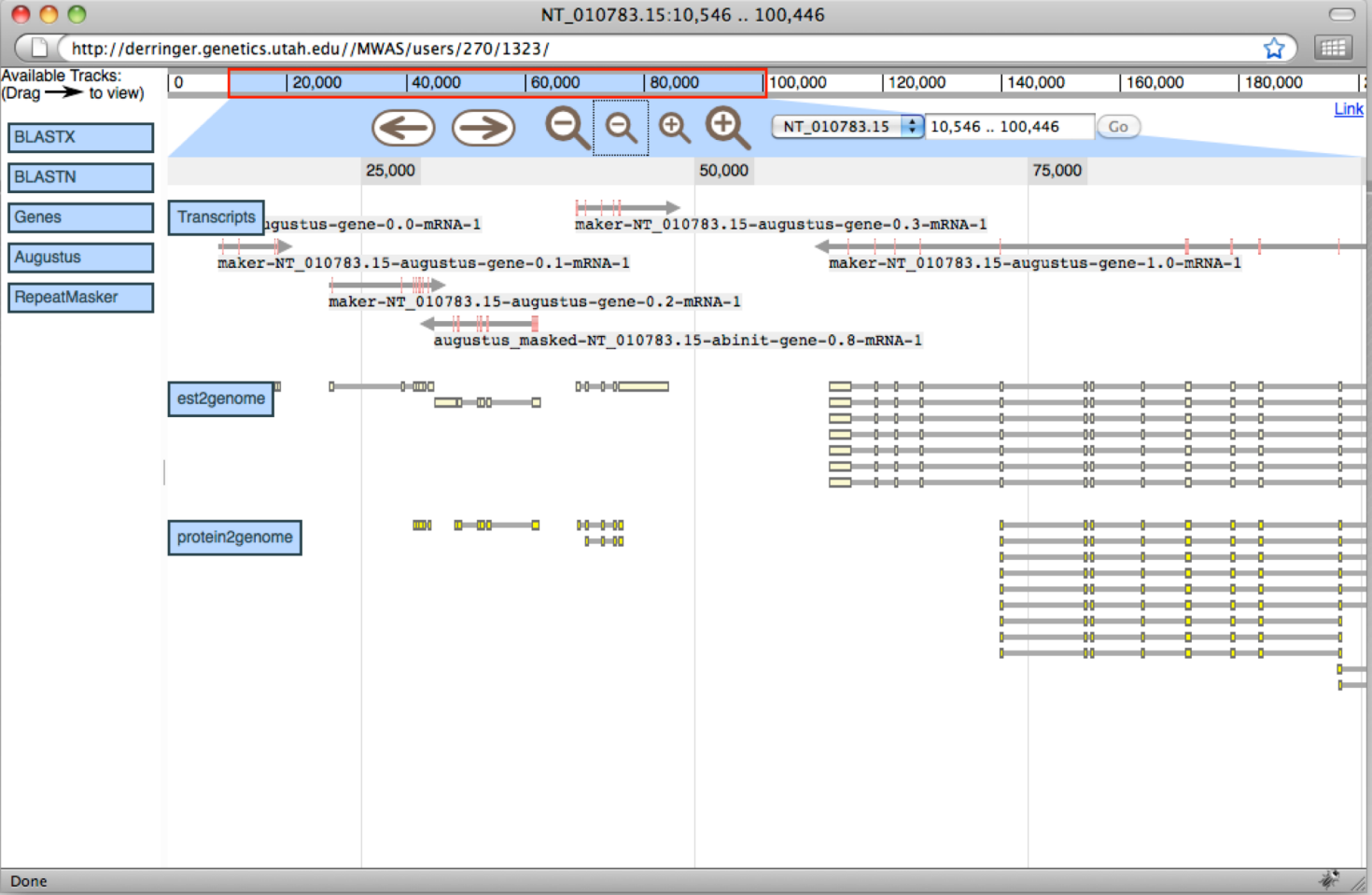
MAKER Gene Annotations

- maker-NT_010783.15-augustus-gene-0.0-mRNA-1
- maker-NT_010783.15-augustus-gene-0.1-mRNA-1
- maker-NT_010783.15-augustus-gene-0.2-mRNA-1
- augustus_masked-NT_010783.15-abinit-gene-0.8-mRNA-1
- maker-NT_010783.15-augustus-gene-0.3-mRNA-1
- maker-NT_010783.15-augustus-gene-1.0-mRNA-1

ESTs - Exonate

- gi|1171543849|ref|NM_173079.2|
- gi|178190461|ref|NM_000988.3|
- gi|1216548122|ref|NM_005533.3|
- gi|118200348|ref|NM_006373.3|
- gi|63252871|ref|NM_007294.2|
- gi|63252872|ref|NM_007295.2|

Done



NT_010783.15:1-201444

File Edit View Tiers Analysis Bookmarks Annotation Window Links Help

0Mb 0.025Mb 0.05Mb 0.075Mb 0.1Mb 0.125Mb 0.15Mb 0.175Mb 0.2Mb

augustus_masked-NT_010783.15-abinit-gene-0.8-mRNA-1

maker-NT_010783.15-augustus-gene-1.0-mRNA-1

Position: [0.103394]

Zoom: [x10] [x2] [x5] [x1] [Reset] Zoom factor = 1.0000

Type	Name	Range	Score
gene	maker-NT_01078...	140963-60144	0.0

maker-NT_010783.15-augustus-gene-1.0			
maker-NT_010783.15-augustus-gene-1.0-mRNA-1			
Genomic Range	Genomic Length	_AED	_QI
61651-60144	1508	0.04	397 0.94 1 1 0.1...
63552-63492	61	0.04	397 0.94 1 1 0.1...
65043-64970	74	0.04	397 0.94 1 1 0.1...
66966-66912	55	0.04	397 0.94 1 1 0.1...

Position: [103394] Feature: [maker-NT_010783.15-augustus-gene-1.0-mRNA-1] Action: []



SOBA

http://www.sequenceontology.org/cgi-bin/soba.cgi/?rm=reload_files&gff_file=NT_010783.15.gff



Home Browser Wiki GFF3 GVF Resources About Request A Term Site Map

Genome Summary
For file(s): NT_010783.15.gff

Below are the feature types and sources for the GFF3 file(s) you uploaded. You must select at least one feature type and source and then click on the headers below to view the analyses.

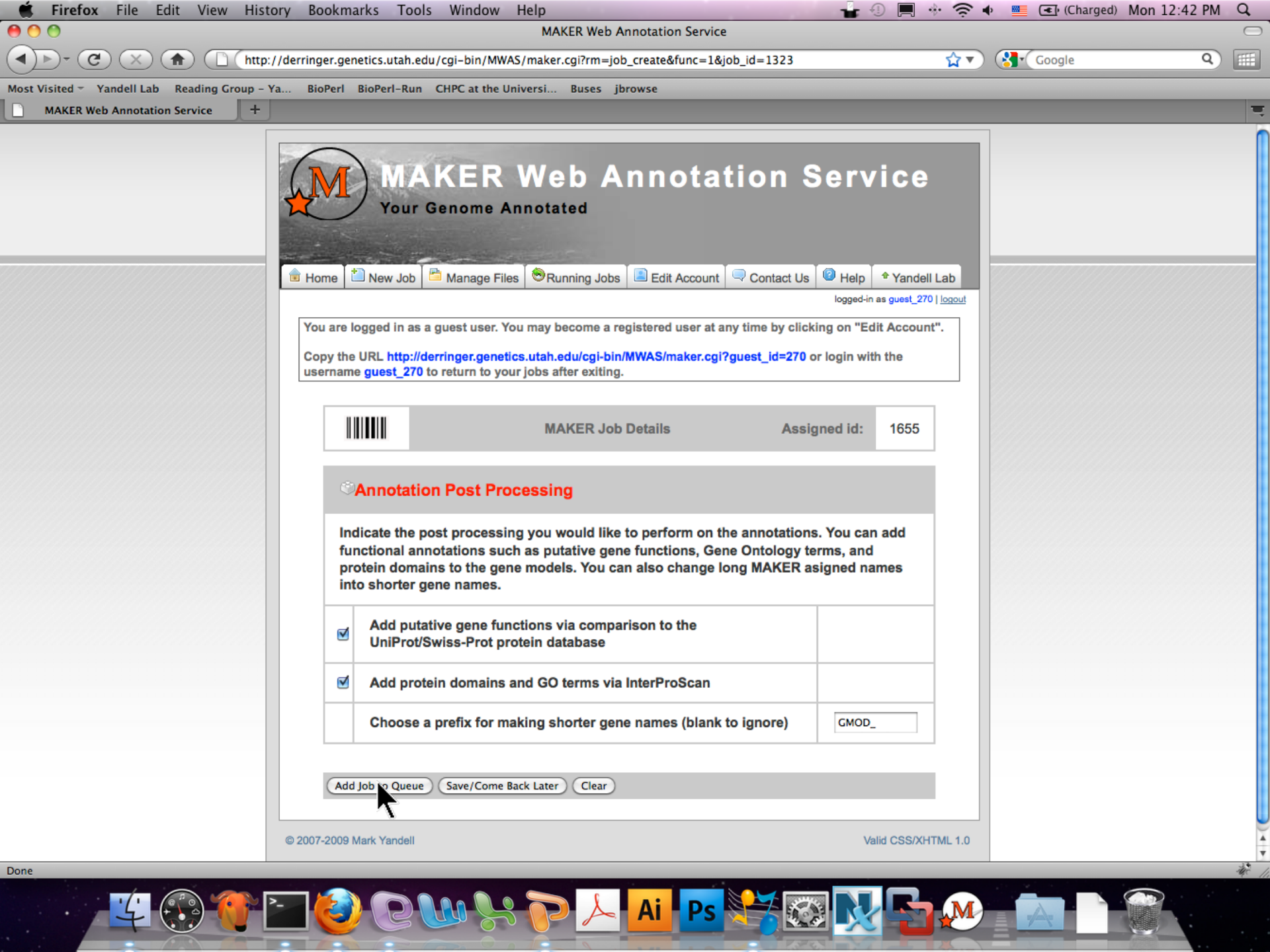
Feature Types			Sources		
Select All	Unselect All	Invert Selection	Select All	Unselect All	Invert Selection
<input checked="" type="checkbox"/> CDS			<input checked="" type="checkbox"/>		
<input checked="" type="checkbox"/> contig			<input checked="" type="checkbox"/> augustus_masked		
<input checked="" type="checkbox"/> exon			<input checked="" type="checkbox"/> blastn		
<input checked="" type="checkbox"/> expressed_sequence_match			<input checked="" type="checkbox"/> blastx		
<input checked="" type="checkbox"/> gene			<input checked="" type="checkbox"/> est2genome		
<input checked="" type="checkbox"/> match			<input checked="" type="checkbox"/> maker		
<input checked="" type="checkbox"/> match_part			<input checked="" type="checkbox"/> protein2genome		
<input checked="" type="checkbox"/> mRNA			<input checked="" type="checkbox"/> repeatmasker		
<input checked="" type="checkbox"/> protein_match					
Select All	Unselect All	Invert Selection	Select All	Unselect All	Invert Selection

Term Validation:
Term usage is OK!

Done

Done





MAKER Web Annotation Service

Your Genome Annotated

- Home
- New Job
- Manage Files
- Running Jobs
- Edit Account
- Contact Us
- Help
- Yandell Lab

logged-in as [guest_270](#) | [logout](#)

You are logged in as a guest user. You may become a registered user at any time by clicking on "Edit Account".
Copy the URL http://derringer.genetics.utah.edu/cgi-bin/MWAS/maker.cgi?guest_id=270 or login with the username [guest_270](#) to return to your jobs after exiting.

 **MAKER Job Details** Assigned id: 1655

Annotation Post Processing

Indicate the post processing you would like to perform on the annotations. You can add functional annotations such as putative gene functions, Gene Ontology terms, and protein domains to the gene models. You can also change long MAKER assigned names into shorter gene names.

<input checked="" type="checkbox"/>	Add putative gene functions via comparison to the UniProt/Swiss-Prot protein database	
<input checked="" type="checkbox"/>	Add protein domains and GO terms via InterProScan	
	Choose a prefix for making shorter gene names (blank to ignore)	<input type="text" value="GMOD_"/>

- Add Job to Queue
- Save/Come Back Later
- Clear



MAKER Web Annotation Service

MAKER Web Annotation Service: 201.4 kbp from NT_010783.15:1..201,444

Overview

Region

Details

NT_010783.15: 70 kbp

MAKER Gene Annotations

- GMOD_00000003-RA: Similar to IFI35: Interferon-induced 35 kDa protein (Homo sapiens)
- GMOD_00000004-RA: Similar to VAT1: Synaptic vesicle membrane protein VAT-1 homolog (Homo sapiens)
- GMOD_00000005-RA: Similar to RND2: Rho-related GTP-binding protein RhoN (Homo sapiens)
- GMOD_00000006-RA: Similar to BRCA1: Breast cancer type 1 susceptibility protein

InterPro Protein Domains

Showing 5 of 22 features

- IPR009909 Nmi/IFP 35
- IPR003579 Ras small GTPase, Rab type
- IPR001357 BRCT
- IPR020843 Polyketide synthase, enoylreductase
- IPR013154 Alcohol dehydrogenase GroES-like

Proteins - Exonerate

- gi|24307901|ref|NP_005524.1
- gi|18379349|ref|NP_006364.2
- gi|4885581|ref|NP_005431.1
- gi|4885069|ref|NP_005159.1
- gi|6552315|ref|NP_009233.1
- gi|63252882|ref|NP_009235.2
- gi|6552301|ref|NP_009226.1
- gi|6552317|ref|NP_009234.1
- gi|6552299|ref|NP_009225.1
- gi|6552321|ref|NP_009236.1
- gi|6552307|ref|NP_009229.1
- gi|6552305|ref|NP_009228.1



GMOD_00000004-RA Details

Name: GMOD_00000004-RA

Type: mRNA

Description: Similar to VAT1: Synaptic vesicle membrane protein VAT-1 homolog (Homo sapiens)

Source: maker

Position: [NT_010783.15:30454..38291 \(- strand\)](#)

Length: 7838

Alias: augustus_masked-NT_010783.15-abinit-gene-0.8-mRNA-1

Dbxref: Gene3D:G3DSA:3.40.50.720

InterPro:IPR002085

InterPro:IPR002364

InterPro:IPR011032

InterPro:IPR013149

InterPro:IPR013154

InterPro:IPR016040

InterPro:IPR020843

PANTHER:PTHR11695

PANTHER:PTHR11695:SF29

Pfam:PF00107

Pfam:PF08240

Prosite:PS01162

SMART:SM00829

superfamily:SSF50129

superfamily:SSF51735

Note: Similar to VAT1: Synaptic vesicle membrane protein VAT-1 homolog (Homo sapiens)

Ontology_term: GO:0003824

GO:0005488

GO:0008152

GO:0008270

GO:0016491

GO:0055114



MAKER Wiki

<http://weatherby.genetics.utah.edu/MAKER/wiki>

Post Processing of Annotations

Once you've determined where the genes are the next question is what do they do. Also it would be nice to change ugly MAKER assigned gene names to follow more standardized formats.

MAKER has a number of accessory scripts that allow you to do just that. So let's take a look at our last example.

```
cd ~/maker_tutorial/example_04_postannotation
ls -l

hsap_contig.gff
hsap_contig.maker.proteins.fasta
hsap_contig.maker.transcripts.fasta
output.blastp
output.iprscan
uniprot_sprot.db
```

Here we have our MAKER output GFF3 and FASTA files for proteins and transcripts ([Click to see GFF3 in JBrowse](#)). We also have output reports from the program [InterProScan](#) (a domain finder) ran on the MAKER proteins and a BLAST report of homology of the MAKER proteins to UniProt/Swiss-Prot. How to run these programs is not part of this tutorial, but how to integrate their output is.

This is an example command line for running BLASTP against UniProt/Swiss-Prot (you don't need to run it, it's just for reference):

```
blastp -query hsap_contig.maker.proteins.fasta -db uniprot_sprot.fasta -evaluate 1e-6 -max_hsps 1 -max_target_seqs 1 -outfmt 6
```

This is an example command line for running InterProScan (you don't need to run it, it's just for reference):

```
interproscan.sh -appl pfam -dp -f TSV -goterms -iprlookup -pa -t p -i hsap_contig.maker.proteins.fasta -o output.iprscan
```

But first lets fix those ugly MAKER names.

MAKER comes with the script `maker_map_ids` to make it easy to rename genes and follow formats such as those suggested by NCBI (organism prefix and gene numbers).

Synopsis:

```
maker_map_ids --prefix PYU1_ --justify 8 genome.all.gff > genome.all.id.map
```

Description:

Acknowledgements

- Mark Yandell
- Barry Moore
- Michael Campbell
- Daniel Ence