

Genome Annotation of *Arabidopsis thaliana* T2T Col-CC By NCBI RefSeq

Terence D. Murphy – RefSeq Eukaryotes and CGR Project Team Lead



National Library of Medicine
National Center for Biotechnology Information

Outline

- The CGR project and NCBI tools for Arabidopsis
- NCBI annotation of Col-CC
- What's Next



NCBI Genome Resources Workshop

Monday January 15, 2024, 4:00 – 6:10 pm, Palm 8

Time	Topic
4:00 – 4:20	SRA: Submission Improvements and Accelerating Discoveries through Lighter Data <i>Yuriy Skripchenko</i>
4:20 – 4:40	Lean, Mean, Genome Cleaning Machine: Removing Contamination with FCS-GX <i>Eric Tvedte</i>
4:40 – 5:00	GenBank Submission: How Do I Submit Diploid Genome Assemblies to NCBI? <i>Shelby Bidwell</i>
5:00 – 5:20	Advances in Eukaryotic Annotation at NCBI <i>Vamsi Kodali</i>
5:20 – 5:40	Using NCBI Datasets to Easily Gather Data from across NCBI Databases <i>Sally Chang</i>
5:40 – 6:00	The NIH Comparative Genomics Resource (CGR) <i>Terence Murphy</i>

Visit NCBI Booth **615**

Contact us info@ncbi.nlm.nih.gov

Watch NCBI News for updates!

<http://www.ncbi.nlm.nih.gov/news/>

<https://www.youtube.com/user/NCBINLM>





Search NCBI ...

[CGR home](#)

[About CGR](#)

[CGR Impact](#)

[Data resources](#)

[Analysis tools](#)

[Data quality tools](#)

[FAQs](#)

NIH Comparative Genomics Resource (CGR)

Maximizing the impact of eukaryotic research organisms

Unlock the full potential of eukaryotic research organisms and their genomic data with the National Institutes of Health (NIH) Comparative Genomics Resource (CGR). CGR facilitates this through community collaboration and an NCBI Toolkit of interconnected and interoperable data and tools.

Get CGR updates

Join our mailing list to receive CGR updates and opportunities to provide feedback.

[Subscribe](#)


https://www.ncbi.nlm.nih.gov/genbank/assembly/col-cc

Genome assembly Col-CC

[Download](#) [datasets](#) [curl](#) Actions

Submitted GenBank assembly	GCA_028009825.2	⋮
Taxon	Arabidopsis thaliana (thale cress)	
Assembly type	haploid	
Submitter	Community-Consensus Arabidopsis Thaliana Reference Genome Assembly Consortium	
Date	Oct 18, 2023	

View the [legacy Assembly page](#)

 [BLAST the reference genome](#)

Assembly statistics

	GenBank
Genome size	142.5 Mb
Total ungapped length	142.5 Mb
Number of chromosomes	5
Number of scaffolds	5
Scaffold N50	27.8 Mb
Scaffold L50	3
Number of contigs	5
Contig N50	27.8 Mb
Contig L50	3
GC percent	37.5
Genome coverage	1000.0x
Assembly level	Complete Genome

assets

Rows per page: 1-20 of 92 < >

Size (Mb)	Level	Release ...	W	Action
119.1	Chromosome	Mar, 2018		⋮
143.5	Complete	Feb, 2023		⋮
142.5	Complete	Oct, 2023		⋮
131.6	Complete	Apr, 2022		⋮
139.9	Chromosome	Feb, 2023		⋮
133.4	Chromosome	Feb, 2023		⋮
133.9	Chromosome	Jul, 2022		⋮
140.0	Chromosome	Feb, 2023		⋮
138.9	Chromosome	Feb, 2023		⋮
148.3	Chromosome	Feb, 2023		⋮
141.0	Chromosome	Feb, 2023		⋮
136.6	Chromosome	Feb, 2023		⋮
139.1	Chromosome	Feb, 2023		⋮
140.4	Chromosome	Feb, 2023		⋮

Comparative Genome Viewer (CGV)

- Compare genomes through assembly alignments
- Over 700 alignments for over 300 species!
- Preprint: PMC10705539

<https://www.ncbi.nlm.nih.gov/genome/cgv/>

The screenshot shows the NCBI Comparative Genome Viewer (CGV) interface. At the top, there is a blue header with the NIH logo and the text "National Library of Medicine National Center for Biotechnology Information". A user profile box in the top right corner displays the name "murphyte2". Below the header is a navigation bar with links for "CGV Home", "Help", and "Release Notes". The main content area is titled "Comparative Genome Viewer" and includes a brief description: "This tool allows you to compare two genomes based on assembly-assembly alignments provided by NCBI." Underneath, there is a section titled "Set up your view" with the instruction: "Make a selection in each of these four steps to view assembly comparison." The interface contains four dropdown menus arranged in a 2x2 grid. The first dropdown is labeled "1. Select a species" and contains "Arabidopsis thaliana (thale cress)". The second dropdown is labeled "2. Select a second species" and contains "Thlaspi arvense". The third dropdown is labeled "3. Select an assembly" and contains "TAIR10.1 (GCF_000001735.4)". The fourth dropdown is labeled "4. Select a second assembly" and contains "T_arvense_v2 (GCA_911865555.2)". Below the dropdowns are two buttons: "Clear Form" and "View Comparison". At the bottom of the form, there is a link: "Not finding your alignment of interest? [Fill out the form](#) to request more alignments." On the right side of the page, there is a yellow "Feedback" button.

CGV: TAIR10 vs Col-CC

First Pass	Total
Col-CC (Current) Coverage: 88.74%	Col-CC (Current) Coverage: 89.58%
TAIR10.1 (Previous) Coverage: 99.03%	TAIR10.1 (Previous) Coverage: 99.32%
Percent Identity: 98.64%	Percent Identity: 98.61%

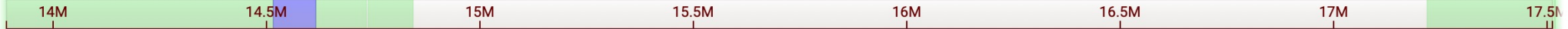
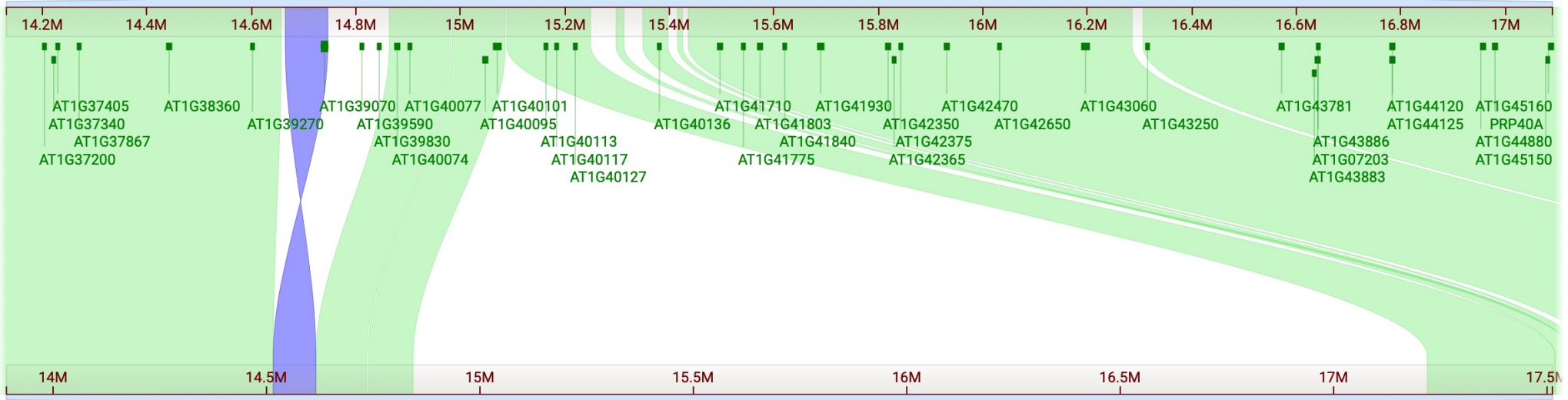
Arabidopsis thaliana Col-CC ([GCA_028009825.2](#))

Arabidopsis thaliana TAIR10.1 ([GCF_000001735.4](#))



Arabidopsis thaliana TAIR10.1 ([GCF_000001735.4](#))

Chr 1



Chr 1

Arabidopsis thaliana Col-CC ([GCA_028009825.1](#))



Outline

- The CGR project and NCBI tools for Arabidopsis
- NCBI annotation of Col-CC
- What's Next

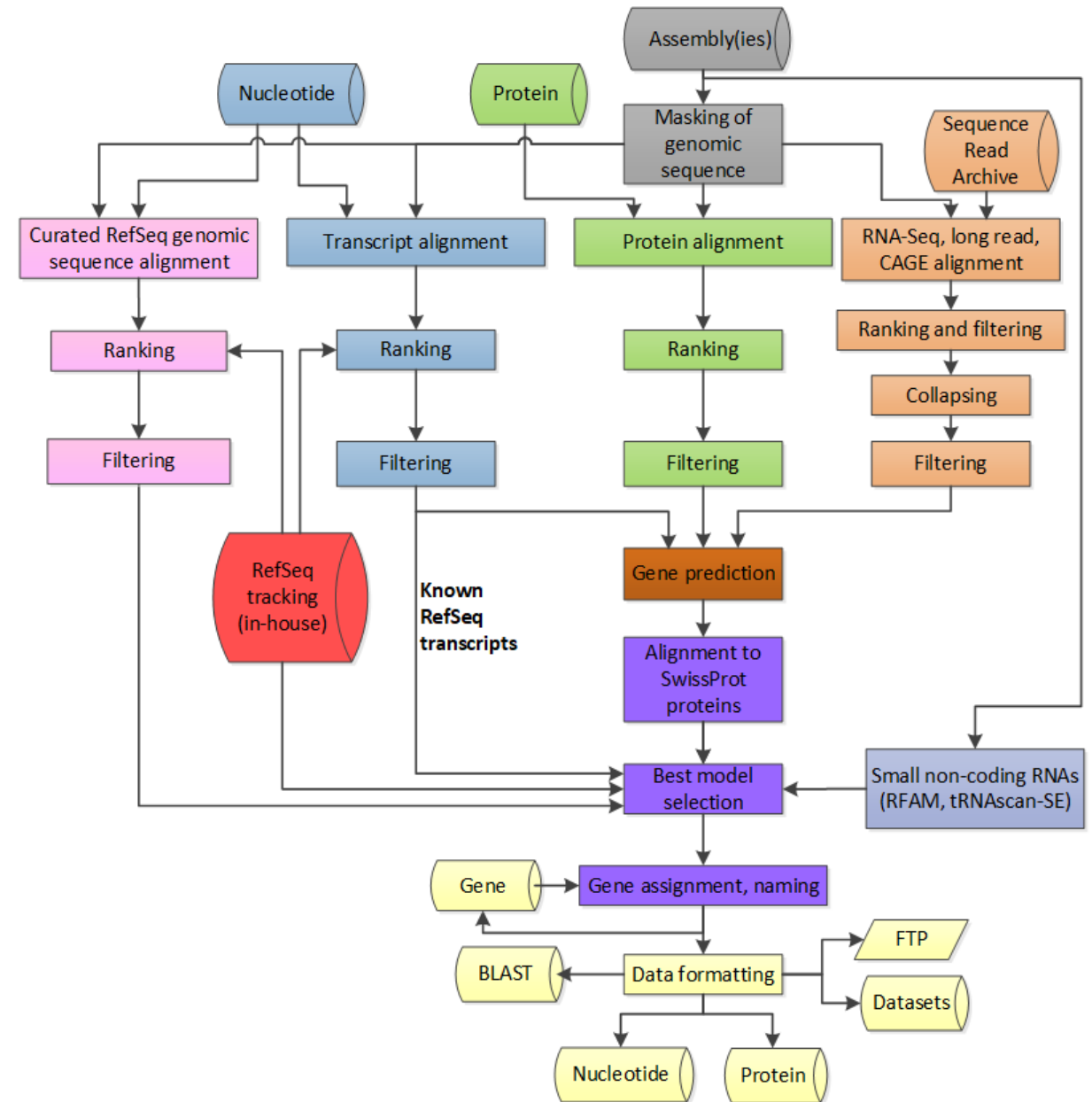


EGAP: Eukaryotic Genome Annotation Pipeline

Used by NCBI to annotate ~1100 species

Evidence used for gene prediction:

- ✓ ESTs
- ✓ cDNAs
- ✓ Same and cross-species proteins
- ✓ RNA-Seq
- ✓ PacBio IsoSeq, ONT transcriptomes
- ✓ CAGE
- ✓ PhyloCSF (for targeted use)



NCBI EGAP annotation

- Software: EGAP v10.1
- Annotation Name: GCA_028009825.1-RS_2023_03
(internal release)

Input type	amount	notes
Oxford Nanopore	41 runs, 61M reads	~90% identity/coverage
PacBio	33 runs, 17M reads	99+% identity/coverage
Illumina	332 runs, 9.3B reads	Mostly unpaired, unstranded
CAGE	52 runs, 3.4M reads	63% aligned
Transcripts	195k cDNA + 1.5M ESTs	99+% identity/coverage
Proteins	<ul style="list-style-type: none">• TAIR A. thaliana annotation• Subset of Brassicaceae proteins	

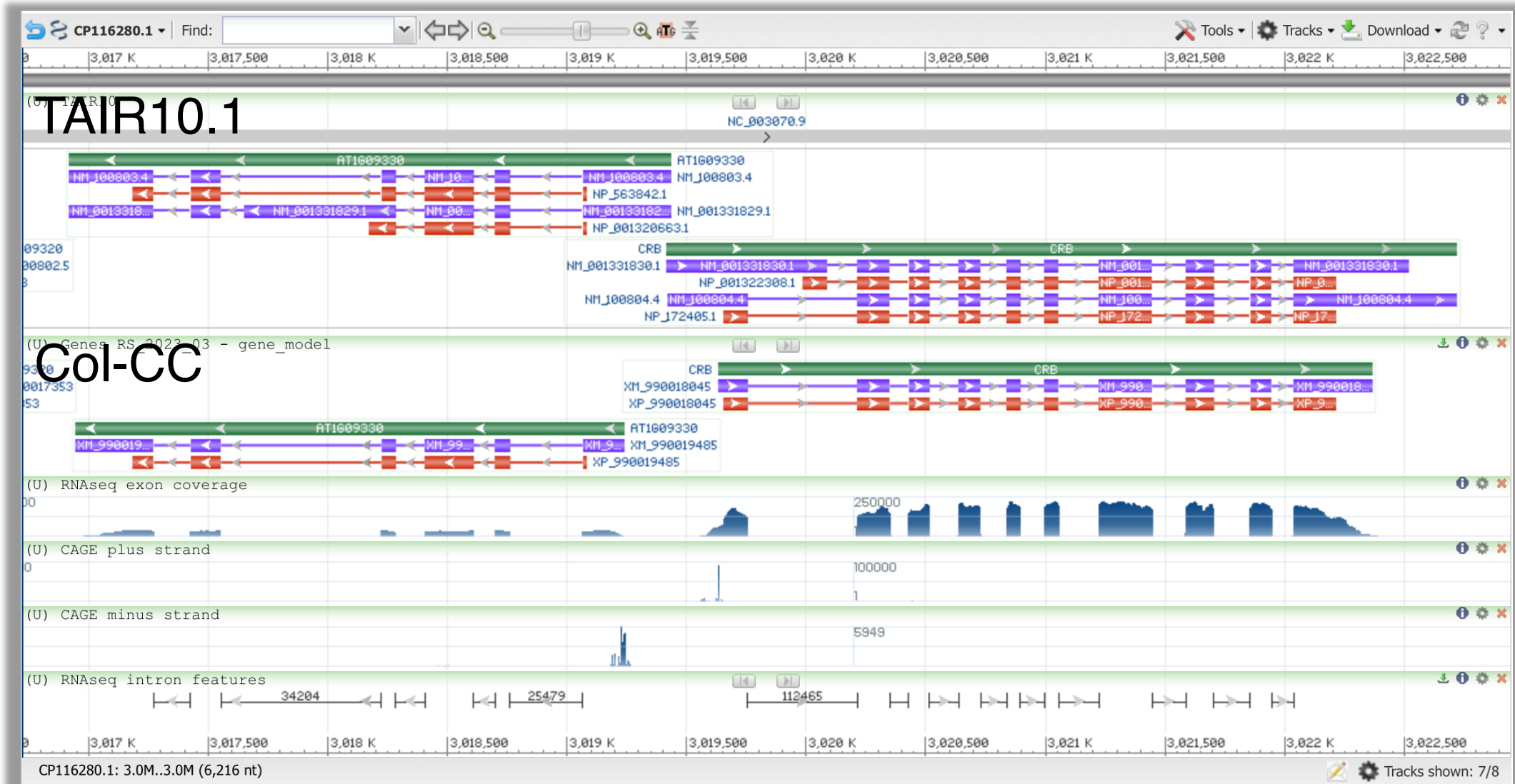
Gene counts

Gene biotype	source	Col-CC	TAIR10.1
Protein-coding	Gnomon	25,565	27,444
with alt transcripts		5,574	10,701
with major corrections		570	
Pseudogene	Gnomon	2,337	4,842
lncRNA	Gnomon	833	3,480
rRNA (v1 assembly)	Rfam	3441	4
snRNA	Rfam	62	82
snoRNA	Rfam	578	287
tRNA	tRNA-scan	545	625
miRNA	previous	325	325
Antisense	previous	79	79

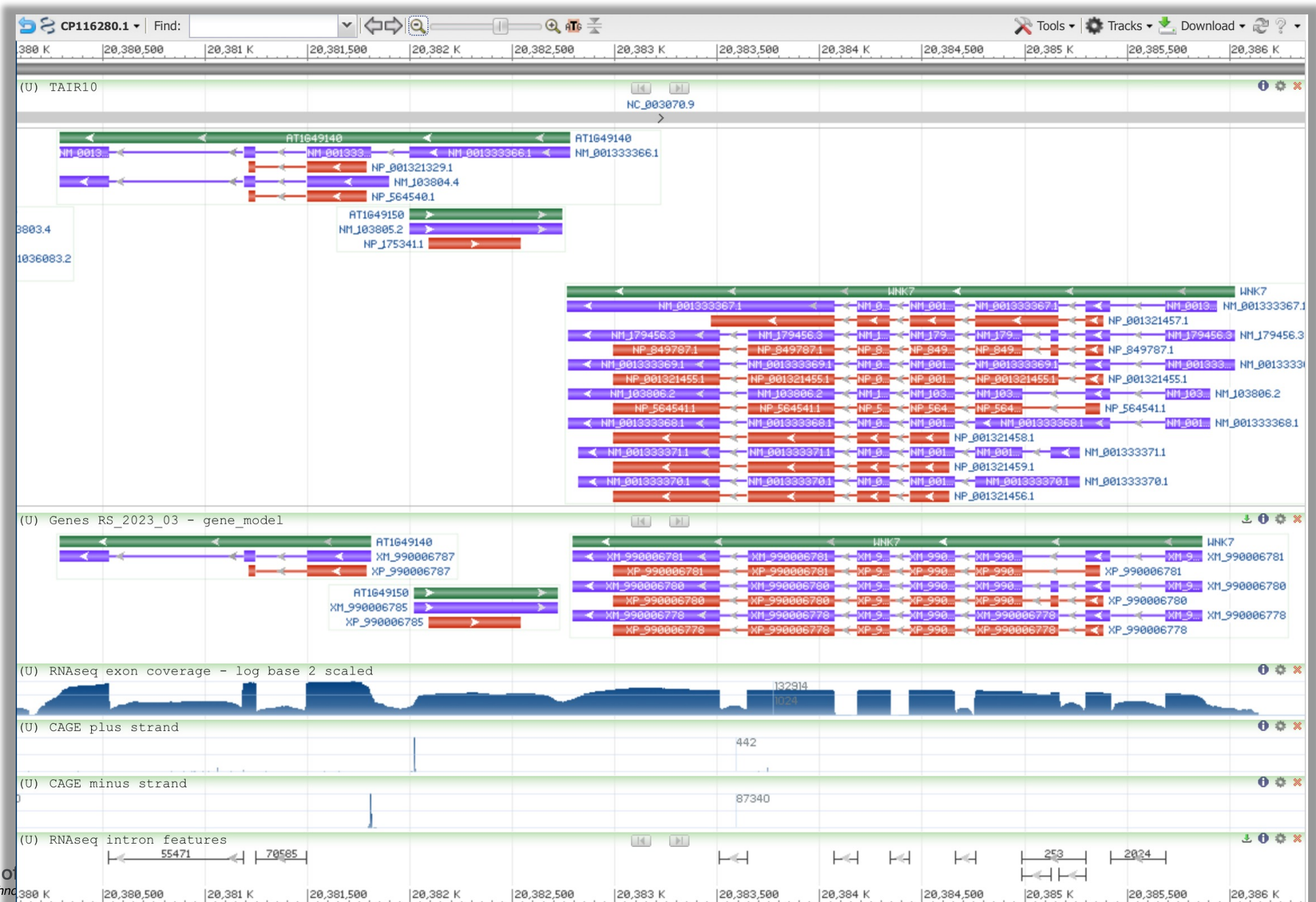
BUSCO:

- TAIR10.1: C:99.8%(S:98.7%,D:1.1%),F:0%,M:0.2%
- RS_2023_03: C:98.8%(S:97.7%,D:1.2%),F:0%,**M1.1%**

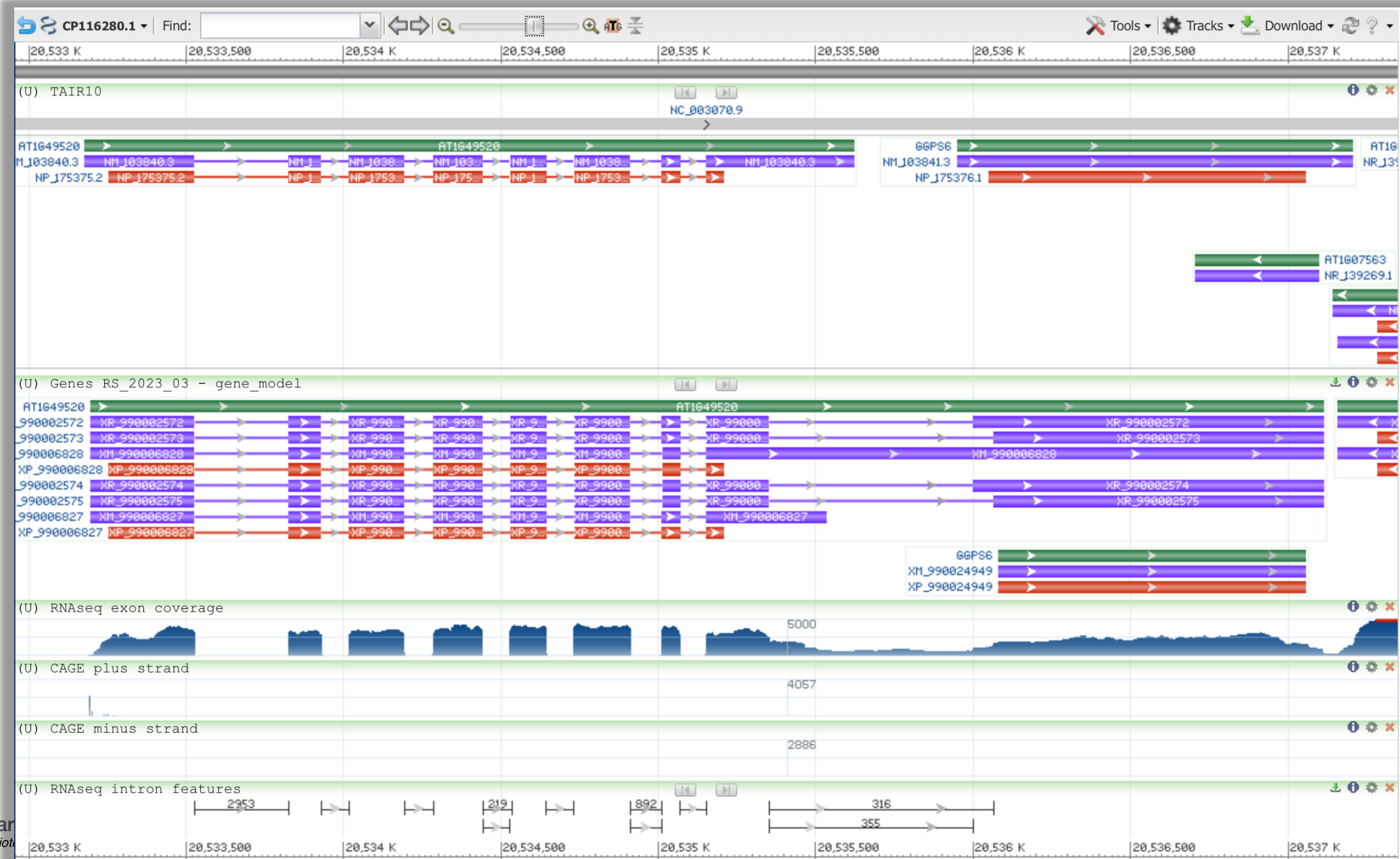
Chr1: AT1G09330 and CRB



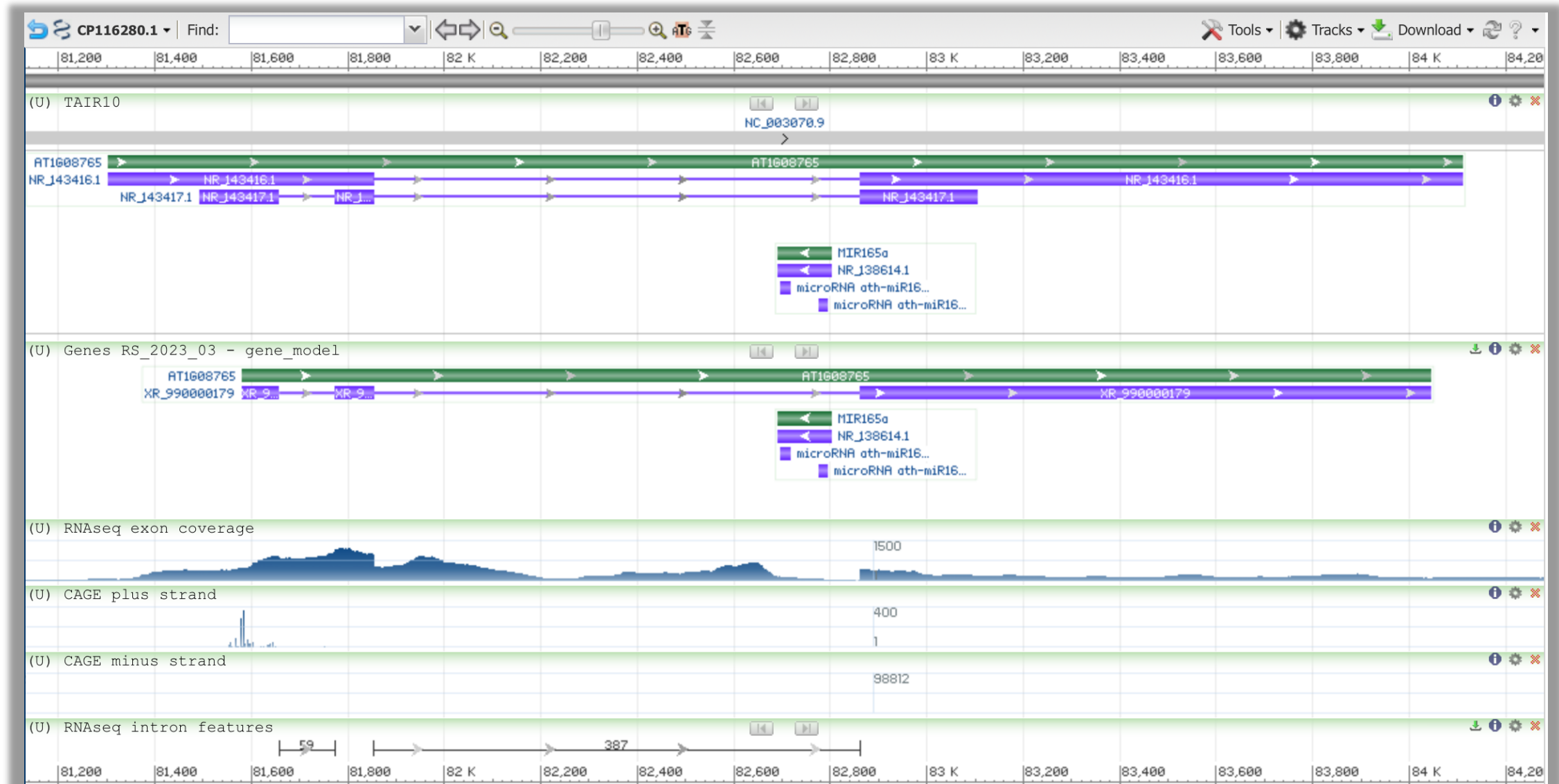
WINK7



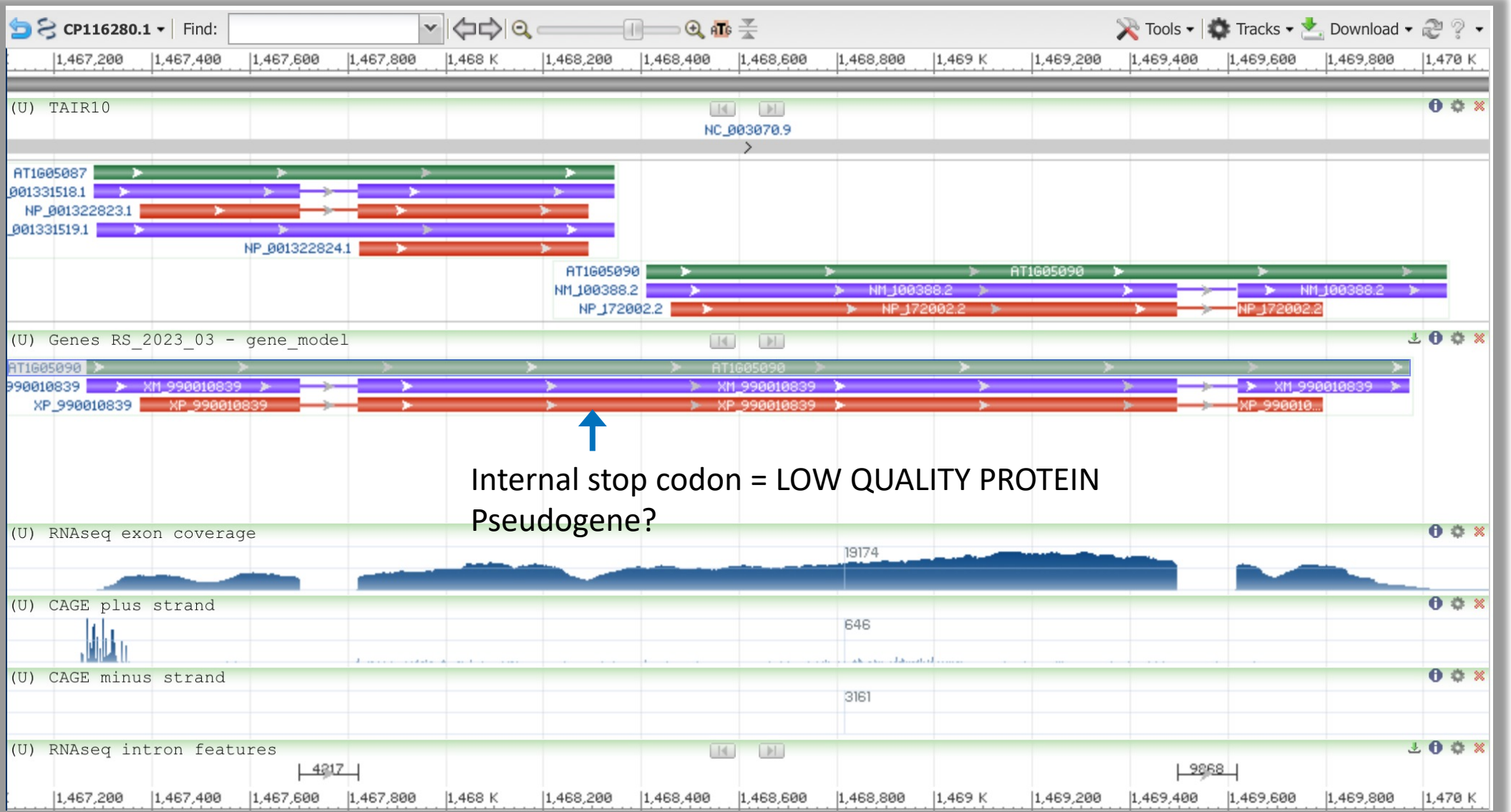
AT1G49520 and GGPS6



AT1G08765 and MIR165a



AT1G05087 and AT1G05090



What's Next



- Review and refinement
 - Major differences vs TAIR10.1
 - Low quality proteins
- Merge other annotations and shift to adjust for added NORs
- Validation and GenBank submission
- Adoption:
 - Incorporate into NCBI RefSeq
 - Switch to Col-CC in TAIR (TAIR11?)
 - Add additional data (expression, variation, etc) to the new reference
- Get back to the real science!

Thank you

This research was supported by the Intramural Research Program of the NIH, National Library of Medicine.

RefSeq/Gene

Shashi Pujar
Tamara Goldfarb
Diana Haddad
John Jackson
Vinita Joardar
Kelly McGarvey
Michael Murphy
Barbara Robbertse
Brian Smith-White
Pooja Strobe
Anjana Vatsan
David Webb

EGAP

Francoise Thibaud-Nissen
Wratko Hlavina
Avi Kimchi
Jinna Hoffman
Vamsi Kodali
Patrick Masterson
Eyal Mozes
Robert Smith
Alexandre Souvorov
Olga Ermolaeva
Craig Wallin

EGAPx / FCS

Alex Astashyn
Nathan Bouk
Victor Joukov
Deacon Sweeney
Eric Tvedte
Victor Sapojnikov
Pooja Strobe
Pape Sylla
Lukas Wagner

CGR

Nuala O'Leary
Sanjida Rangwala
Tom Madden
Aron Marchler-Bauer
Anne Ketter
Katya Sukharnikov

NCBI Leadership

Valerie Schnieder
Kim Pruitt
Steve Sherry

Visit NCBI Booth **615**

Contact us info@ncbi.nlm.nih.gov

Watch NCBI News for updates!

<http://www.ncbi.nlm.nih.gov/news/>

<https://www.youtube.com/user/NCBINLM>

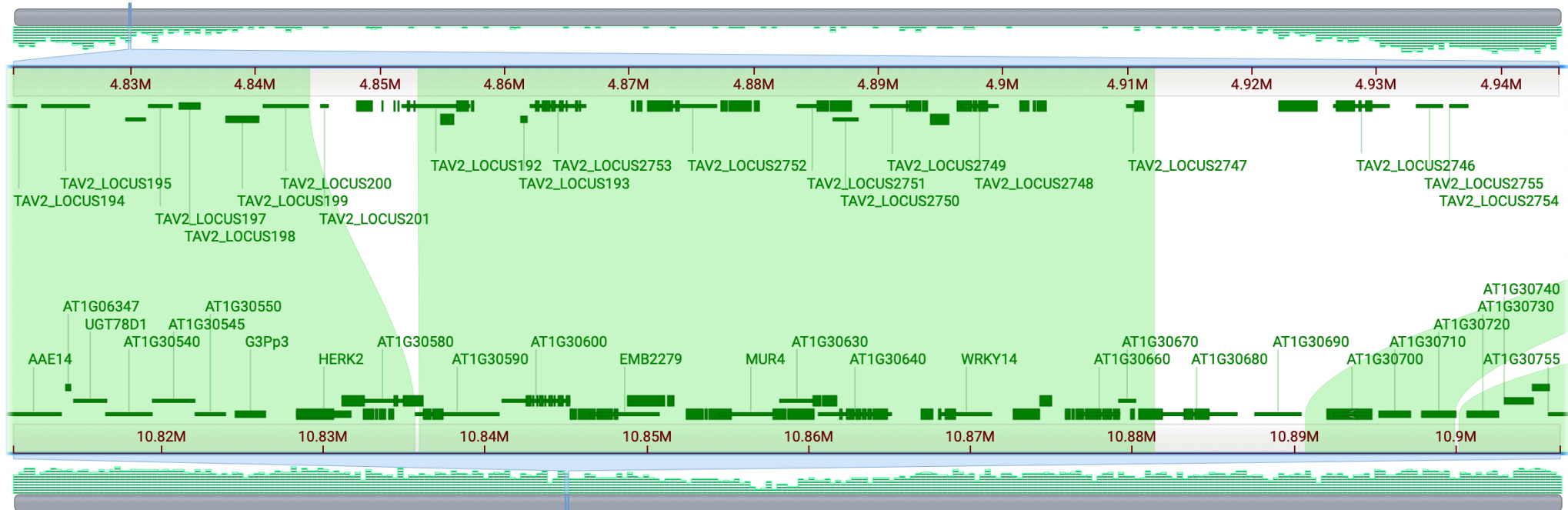
CGV: Arabidopsis vs Thlaspi

Thlaspi arvense T_arvense_v2 ([GCA_911865555.2](#))

Thlaspi arvense T_arvense_v2 ([GCA_911865555.2](#))

Thlaspi arvense T_arvense_v2 ([GCA_911865555.2](#))

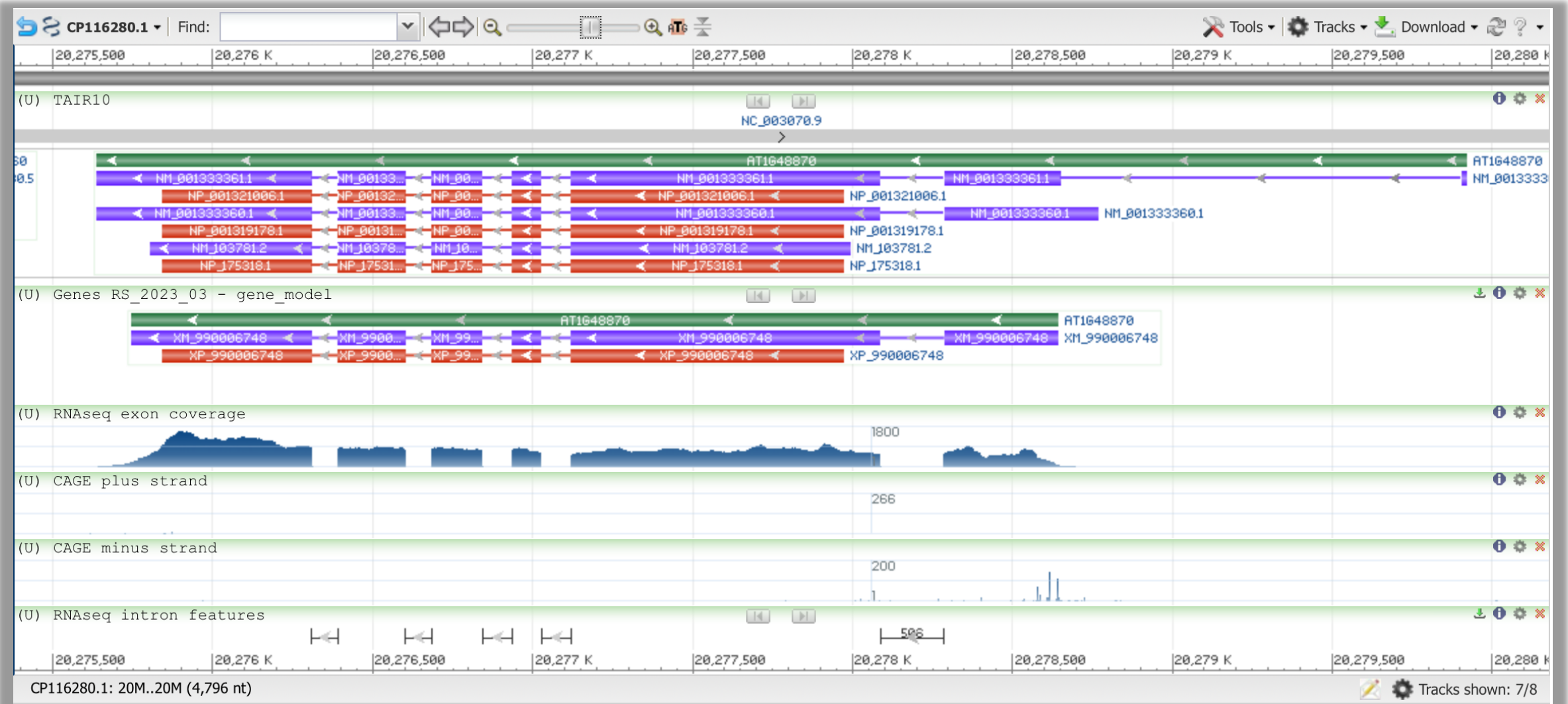
Chr 1



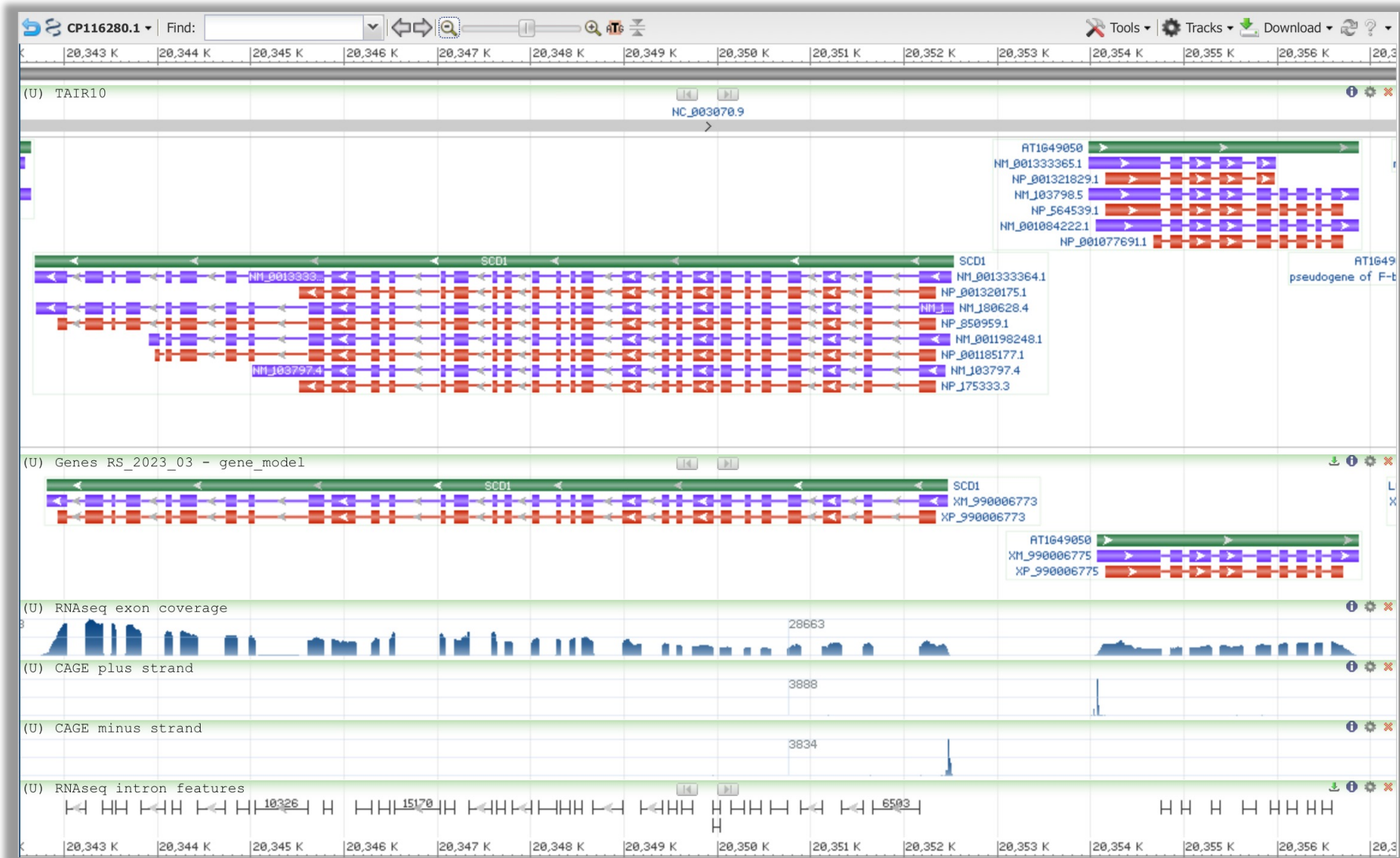
Chr 1

Arabidopsis thaliana TAIR10.1 ([GCF_000001735.4](#))

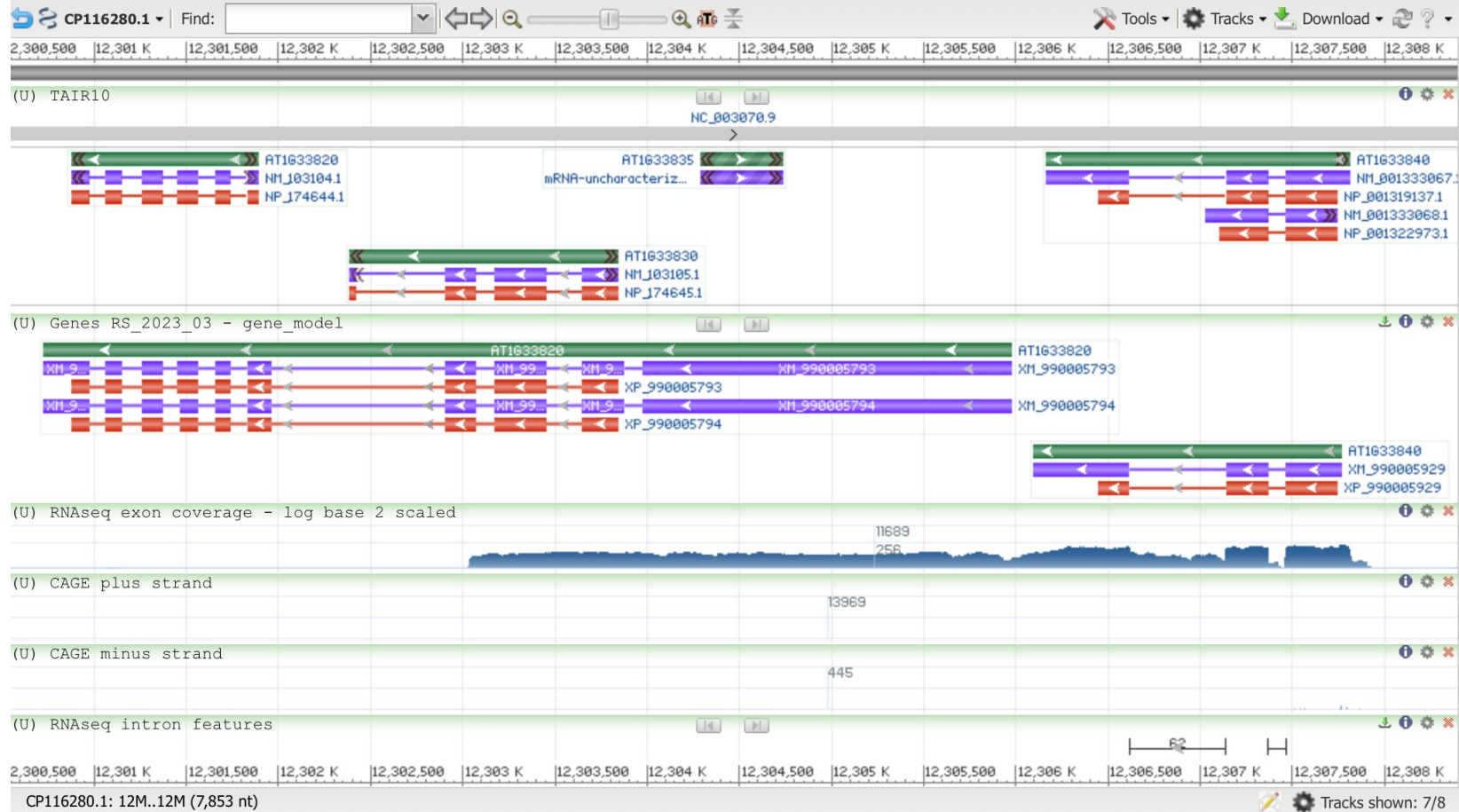
Chr1: AT1G48870



Chr1: SCD1 and AT1G49050



AT1G33820 and AT1G33840



AT1G36920 and AT1G36925

