



Evaluation of Gene Finding Tools on *Trichoderma* Genomes - Update

July, 2025

Connor Burbridge



Outline

- Progress update with results
- Progress assessment
- Timeline for completion
- Discussion of next steps



Background

- Gene finding is a critical step in genome annotation.
- *Trichoderma* species are important for agriculture and biotechnology.
- Various tools exist for gene prediction, but their implementations and performance can vary significantly.
- Few studies have compared these tools in fungi, and even fewer in *Trichoderma*.
- In addition, high-quality reference genomes are not available for all species, which makes comparisons difficult.



Why *Trichoderma*?

- *Trichoderma* species are ubiquitous in soil and play a significant role in nutrient cycling.
- They are also used in biocontrol and as biofertilizers.
- Known for their production of plant cell wall degrading enzymes, and secondary metabolites.
- Further genomic studies can help in understanding their biology and potential applications.

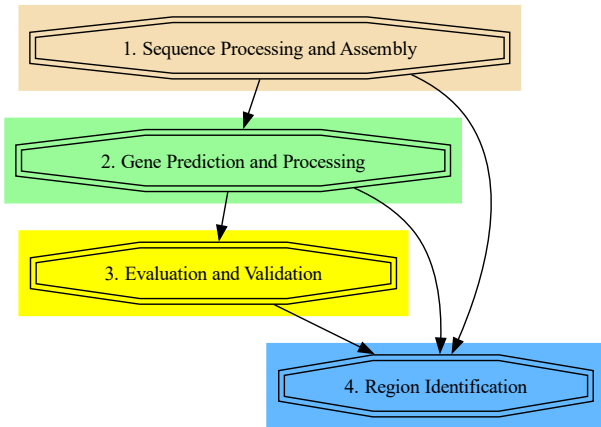


Novel *Trichoderma* Genomes

- *Trichoderma* species are diverse, with many species not yet fully characterized.
- Recent advances in sequencing technology have made it possible to generate high-quality genomes for these species.
- Two novel genomes of *Trichoderma* species have been sequenced at the Global Institute for food Security
- DC1 and Tsth20, shown to improve drought and salt tolerance when applied to crops.
- These genomes provide an opportunity to evaluate gene finding tools in a comparative context, and contribute to the understanding of *Trichoderma* biology.



Workflow Overview



Datasets and Tools

- **Datasets:**

- *Trichoderma* genomes: DC1 and Tsth20.
- Reference genomes and annotations: *T. reesei*, *T. harzianum*, *T. virens*.
- RNAseq training data from *T. reesei*.
- Benchmarking Universal Single-Copy Orthologs (BUSCO) fungal database.
- Protein sequence queries for tblastn from *T. atroviride*, *Fusarium graminearum*, and *Saccharomyces cerevisiae*.

- **Tools:**

- Sequence processing: FastQC, Trimmomatic, Hisat2.
- Genome assembly: NextDenovo and NextPolish.
- Gene finding tools: Braker2 and GeneMark-ES.
- Evaluation tools: BUSCO, tblastn, InterProScan, and custom scripts.



Assembly Results

Name	Total Contigs	Total Length	Largest Contig	GC%	N50	L50
DC1	8	38.6 Mb	11.49 Mb	47.97	5.69 Mb	3
Tsth20	7	41.58 Mb	8.02 Mb	47.33	6.52 Mb	3
<i>T. harzianum</i>	532	40.98 Mb	4.08 Mb	47.61	2.41 Mb	7
<i>T. virens</i>	93	39.02 Mb	3.45 Mb	49.25	1.83 Mb	8
<i>T. reesei</i>	77	33.39 Mb	3.75 Mb	52.82	1.21 Mb	9

- DC1 and Tsth20 assemblies are high quality with few contigs in comparison to other *Trichoderma* assemblies.
- Input sequences and assemblies show a bimodal distribution of GC content.
- AT-rich sequence content may be related to transposable elements and repeat-induced mutations, which may be of interest in secondary metabolite production.



Gene Finding Results

Assembly	Braker2		GeneMark		RefSeq	
	Genes	CDS	Genes	CDS	Genes	CDS
DC1	8546	8637	11353	11353	N/A	N/A
Tsth20	8784	8858	12362	12362	N/A	N/A
<i>T. reesei</i>	9659	10175	9196	9196	9109	9118
<i>T. harzianum</i>	8314	8385	12164	12164	14269	14090
<i>T. virens</i>	7801	7863	11866	11866	12405	12406

- Braker2 predicts more genes and coding sequences than GeneMark and RefSeq in *T. reesei*, but fewer in other assemblies.
- GeneMark and RefSeq predictions are similar, except in *T. harzianum*, where RefSeq predicts 17% more genes.
- GeneMark only predicts one coding sequence per gene, while Braker2 and RefSeq predict multiple coding sequences.





Examining Coding Sequence Lengths

Genome	Tool #1	Tool #2	P-value
DC1	Braker2	GeneMark	0.999
Tsth20	Braker2	GeneMark	0.965
<i>T. reesei</i>	Braker2	GeneMark	$9.481 * 10^{-07}$
<i>T. reesei</i>	GeneMark	RefSeq	0.002
<i>T. reesei</i>	Braker2	RefSeq	$1.340 * 10^{-07}$
<i>T. harzianum</i>	Braker2	GeneMark	0.863
<i>T. harzianum</i>	GeneMark	RefSeq	$4.313 * 10^{-52}$
<i>T. harzianum</i>	Braker2	RefSeq	$4.674 * 10^{-55}$
<i>T. virens</i>	Braker2	GeneMark	0.635
<i>T. virens</i>	GeneMark	RefSeq	$7.352 * 10^{-12}$
<i>T. virens</i>	Braker2	RefSeq	$1.794 * 10^{-09}$



BUSCO Results

Strain	Complete	Single	Duplicated	Fragmented	Missing
DC1	99.5	80.2	19.3	0.1	0.4
Tsth20	99.9	81.7	18.2	0.0	0.1
<i>T. harzianum</i>	99.7	80.2	19.5	0.0	0.3
<i>T. virens</i>	99.8	79.0	20.8	0.1	0.1
<i>T. reesei</i>	99.9	85.5	14.4	0.1	0.0

(a) Braker2

Strain	Complete	Single	Duplicated	Fragmented	Missing
DC1	99.2	98.8	0.4	0.3	0.5
Tsth20	99.8	99.1	0.7	0.0	0.2
<i>T. harzianum</i>	99.6	98.9	0.7	0.0	0.4
<i>T. virens</i>	99.7	99.2	0.5	0.1	0.2
<i>T. reesei</i>	99.6	99.5	0.1	0.0	0.4

(b) GeneMark

Strain	Complete	Single	Duplicated	Fragmented	Missing
<i>T. harzianum</i>	99.9	99.2	0.7	0.0	0.1
<i>T. virens</i>	99.5	98.8	0.7	0.3	0.2
<i>T. reesei</i>	99.8	99.5	0.3	0.0	0.2

(c) RefSeq

Agreement of Gene Predictions

- In DC1 and Tsth20, BRaker2 and GeneMark tend to agree on the start and stop positions of genes when they predict the same gene. Both gene finders have a large portion of singleton predictions.

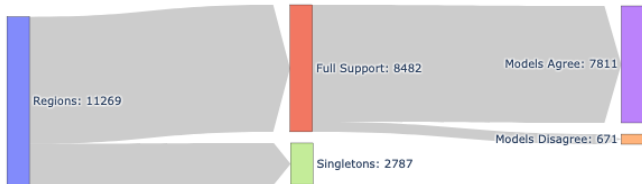


Figure: DC1

- In *T. reesei*, Braker2, GeneMark and RefSeq tend to disagree more on the start and stop positions of genes when they predict the same gene.
- Disagreement is more pronounced in *T. harzianum* and *T. virens*, where genes with partial support from more than one gene finder are common.

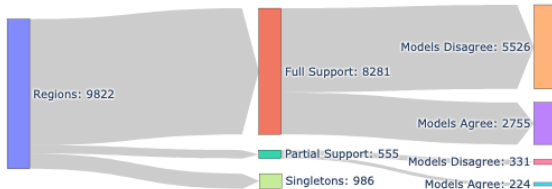


Figure: *T. reesei*

InterProScan Results

- I need to look into this. Braker2 numbers strange

Assembly	Braker2	GeneMark	RefSeq
DC1	<u>10676</u> <u>14479</u>	<u>8416</u> <u>11354</u>	N/A
Tsth20	<u>11389</u> <u>15546</u>	<u>9168</u> <u>12373</u>	N/A
<i>T. reesei</i>	<u>8471</u> <u>11704</u>	<u>6990</u> <u>9196</u>	<u>6964</u> <u>9111</u>
<i>T. harzianum</i>	<u>11370</u> <u>15408</u>	<u>9061</u> <u>12164</u>	<u>9293</u> <u>14065</u>
<i>T. virens</i>	<u>11249</u> <u>15062</u>	<u>8871</u> <u>11866</u>	<u>9062</u> <u>12383</u>



BLAST Results

Reference	Ref. Proteins	DC1	Tsth20	<i>T. reesei</i>	<i>T. harzianum</i>	<i>T. virens</i>
<i>Trichoderma atroviride</i>	11807	11552	11080	10601	11081	11078
<i>Fusarium graminearum</i>	13312	10327	10429	10064	10434	10490
<i>Saccharomyces cerevisiae</i>	6014	3537	3517	3445	3509	3500

- The *T. atroviride* and *Fusarium* datasets are well represented in the tblastn searches, while the *S. cerevisiae* dataset is less well represented.





Subject	Query	Braker2	GeneMark	RefSeq
DC1	<i>T. atroviride</i>	5902	4679	N/A
DC1	<i>F. graminearum</i>	4955	4114	N/A
DC1	<i>S. cerevisiae</i>	2105	1850	N/A
<i>T. reesei</i>	<i>T. atroviride</i>	5072	5174	4989
<i>T. reesei</i>	<i>F. graminearum</i>	4577	4685	4529
<i>T. reesei</i>	<i>S. cerevisiae</i>	2055	2114	2022
<i>T. harzianum</i>	<i>T. atroviride</i>	6363	4611	6835
<i>T. harzianum</i>	<i>F. graminearum</i>	5659	4198	5982
<i>T. harzianum</i>	<i>S. cerevisiae</i>	2424	1963	2560

- Regions with Braker2 and RefSeq gene predictions consistently have more tblastn hits than GeneMark, with the exception of *T. reesei*.



Gene Predictions in AT-rich Sequence

Assembly	Braker2	GeneMark	RefSeq
DC1	31	11	N/A
Tsth20	11	2	N/A
<i>T. reesei</i>	39	48	107
<i>T. harzianum</i>	81	30	154
<i>T. virens</i>	21	8	20

- Very few genes are predicted in AT-rich regions of DC1, Tsth20 and *T. virens* assemblies in comparison to the *T. reesei* and *T. harzianum* assemblies.
- Possibly due to higher quality assemblies? Although *T. virens* is still fragmented.



Tool	DC1	Tsth20	<i>T. reesei</i>	<i>T. harzianum</i>	<i>T. virens</i>
Braker2	$9.56 * 10^{-181}$	$1.14 * 10^{-259}$	$2.68 * 10^{-96}$	$4.05 * 10^{-140}$	$1.35 * 10^{-35}$
GeneMark	$5.12 * 10^{-216}$	0.0	$5.66 * 10^{-49}$	$5.37 * 10^{-219}$	$5.31 * 10^{-35}$
RefSeq	N/A	N/A	$1.29 * 10^{-49}$	$2.44 * 10^{-205}$	$7.40 * 10^{-33}$

- The results of the two-sided binomial test indicate that the number of genes predicted in AT-rich genomic sequence is not proportional to the fraction of genomic sequence they comprise.



Category	Braker2	GeneMark	RefSeq
Availability	3	3	0
Ease of install	1	2	0
Ease of use	3	3	0
<hr/>			
# of genes predicted	0	3	3
# of transcripts predicted	3	0	2
Predicts shortest genes	2	1	0
Predicts more shorter genes	1	0	3
BUSCO Performance	2	1	3
Performance in AT-rich sequence	2	1	3
Predictions with InterProScan support	3	3	3
Final Score (Publicly Available)	20	17	N/A
Final Score (Ignoring Availability)	13	9	17

