

COMPARATIVE ANALYSIS OF GENE FINDING TOOLS WHEN
APPLIED TO *Trichoderma* GENOMES

A thesis submitted to the
College of Graduate and Postdoctoral Studies
in partial fulfillment of the requirements
for the degree of Master of Science
in the Applied Computing of Computer Science
University of Saskatchewan
Saskatoon

By
Connor Burbridge

©Connor Burbridge, All rights reserved.

Unless otherwise noted, copyright of the material in this thesis belongs to
the author.

Permission to Use

In presenting this thesis in partial fulfillment of the requirements for a Postgraduate degree from the University of Saskatchewan, I agree that the Libraries of this University may make it freely available for inspection. I further agree that permission for copying of this thesis in any manner, in whole or in part, for scholarly purposes may be granted by the professor or professors who supervised my thesis work or, in their absence, by the Head of the Department or the Dean of the College in which my thesis work was done. It is understood that any copying or publication or use of this thesis or parts thereof for financial gain shall not be allowed without my written permission. It is also understood that due recognition shall be given to me and to the University of Saskatchewan in any scholarly use which may be made of any material in my thesis.

Disclaimer

Reference in this thesis to any specific commercial products, process, or service by trade name, trademark, manufacturer, or otherwise, does not constitute or imply its endorsement, recommendation, or favoring by the University of Saskatchewan. The views and opinions of the author expressed herein do not state or reflect those of the University of Saskatchewan, and shall not be used for advertising or product endorsement purposes.

Requests for permission to copy or to make other uses of materials in this thesis in whole or part should be addressed to:

Head of the Department of Computer Science
176 Thorvaldson Building, 110 Science Place
University of Saskatchewan
Saskatoon, Saskatchewan S7N 5C9 Canada

OR

Dean
College of Graduate and Postdoctoral Studies
University of Saskatchewan
116 Thorvaldson Building, 110 Science Place
Saskatoon, Saskatchewan S7N 5C9 Canada

Abstract

This thesis presents a comparative analysis of gene prediction tools applied to genomes of the genus *Trichoderma*. The study evaluates the performance of several widely used gene finding algorithms, assessing their accuracy, sensitivity, and specificity using both computational benchmarks and biological validation. Genomic datasets from multiple *Trichoderma* species were analyzed, and the predicted gene models were compared against reference annotations and experimental data. The results highlight the strengths and limitations of each tool, providing recommendations for optimal gene annotation strategies in fungal genomics. Briefly summarizing the results, we found that gene finders generally showed high sensitivity in detecting conserved genes, and predicted gene models contained functional domains, as validated by InterProScan. However, discrepancies were observed in start and stop positions, exon-intron structures, and alternative splicing events among the tools. Notably, gene prediction accuracy declined in AT-rich genomic regions, where fewer genes were identified compared to regions with normal GC content. When selecting a gene finder, users should consider several factors, including available annotations and their quality, the specific genomic context, and the availability of appropriate training data. If a user is working with a representative organism, we recommend they use the annotation available for that organism, but they should consider the quality of the genome assembly and the assembly methods used. If a user is not working with a representative organism, we recommend they use Braker2 only if they have access to training data from the organism of interest or a closely related organism; otherwise, we recommend they use GeneMark. This work contributes to improved annotation pipelines and a deeper understanding of gene structure in *Trichoderma*, facilitating future research in genetics, biotechnology, and agriculture.

In addition to the main findings, the genomes of two novel *Trichoderma* species, DC1 and Tsth20, were successfully assembled and annotated, resulting in high-quality genomic resources. The assemblies exhibited near-chromosomal scale contigs and performed well in standard assembly metrics. The gene predictions generated by the evaluated tools produced a rich dataset of annotations for these genomes, which can be leveraged for further biological investigations.

Acknowledgements

I would like to acknowledge both Dr. Kusalik and Dave Schneider for their excellent support and mentorship during both COVID and my MSc. project. I would also like to thank Brendan Ashby and Dr. Leon Kochian for providing both data and additional financial support. I also thank my family and girlfriend for their unwavering support and encouragement throughout my studies.

Contents

Permission to Use	i
Abstract	ii
Acknowledgements	iii
Contents	iv
List of Tables	vi
List of Figures	vii
List of Abbreviations	viii
1 Introduction	1
2 Background	3
2.1 <i>Trichoderma</i> Species and their Features	3
2.2 Novel <i>Trichoderma</i> Genomes	4
2.3 Evolutionary Landscape of <i>Trichoderma</i>	5
2.4 Gene Predictions and Complementing Similarity Searches	5
2.5 Secondary Metabolites	8
2.6 Gene Prediction Tools in Fungi	9
3 Research Questions	11
3.1 <i>Trichoderma</i> Assembly Results	11
3.2 Profiling of Gene Finding Tools	11
3.3 Number of Features Predicted	11
3.4 Lengths of Predicted Genes	12
3.5 BUSCO Completeness	12
3.6 Validation of Predicted Genes via InterProScan	12
3.7 Similarity Searches of Predicted Genes with tblastn	12
3.8 Identifying Candidate Secondary Metabolite Biosynthetic Gene Clusters	13
3.9 Agreement Between Gene Finders	13
3.10 Performance in AT-rich Genomic Sequence	13
3.11 Selection of a Gene Finding Tool	13
4 Data and Methodology	14
4.1 Methodology Overview	14
4.2 Sequence Processing and Assembly Workflow	14
4.2.1 Sequence Pre-processing and Assembly	16
4.2.2 Identifying AT-rich Genomic Sequence	16
4.3 Overview of Gene Prediction and Further Processing	18
4.3.1 Selection of Datasets from NCBI	19
4.3.2 Gene Prediction Using GeneMark-ES	21
4.3.3 Gene Prediction Using Braker2	21
4.3.4 Statistical Analysis of CDS Lengths	21
4.4 Overview of Evaluation and Validation of Gene Predictions	22
4.4.1 Evaluating Gene-Finder Performance Using BUSCO	22
4.4.2 Functional Annotation Using InterProScan	22

4.4.3	Ground-truthing with the Basic Local Alignment Search Tool (BLAST)	23
4.4.4	Identification of Secondary Metabolite Biosynthetic Gene Clusters Using antiSMASH	23
4.5	Region Identification Overview	24
4.5.1	Region Identification Algorithm	24
4.5.2	Regions in AT-rich Genomic Sequence	26
4.5.3	Visualization of Identified Regions	26
5	Results and Discussion	28
5.1	Overview	28
5.2	Assemblies of DC1 and Tsth20	28
5.3	Installation and Profiling Gene Finding Tools	31
5.4	Initial Gene Finding Results	32
5.5	Distribution of Predicted Coding Sequence Lengths	34
5.6	BUSCO Results	38
5.7	Identifying Secondary Metabolite Biosynthetic Gene Clusters	41
5.8	Region Identification	44
5.9	InterProScan as Supporting Evidence for Predicted Genes	48
5.10	BLAST Results	50
5.11	Genes in Regions of Anomalous GC Content	52
6	Conclusions	60
6.1	Selection of Gene Finding Tool	60
6.2	Future Work	61
	References	64

List of Tables

5.1	General assembly metrics produced by QUAST	29
5.2	Number of genes and coding sequences predicted	33
5.3	Results of Kolmogorov-Smirnov tests	38
5.4	BUSCO results	40
5.5	Missing BUSCO IDs	42
5.6	Additional missing BUSCO IDs	44
5.7	InterProScan Pfam Evidence	50
5.8	Overview of tblastn alignments	51
5.9	Counts of transcripts in regions with tblastn hits	52
5.10	Agreement of gene predictions in AT-rich regions	53
5.11	Number of genes predicted in AT-rich regions	54
5.12	Binomial test results	54
6.1	Final scoring table	62

List of Figures

2.1	<i>Trichoderma</i> Phylogenetic tree	6
2.2	Generic gene model	7
4.1	Workflow overview	15
4.2	Sequence assembly and processing workflow	17
4.3	GC content distribution for DC1	18
4.4	Gene prediction workflow	20
4.5	Gene evaluation and validation workflow	23
4.6	Region identification workflow	25
4.7	IGV example	26
5.1	GC content distribution for each <i>Trichoderma</i> genome assembly	30
5.2	Number of genes predicted	33
5.3	Number of coding sequences predicted	34
5.4	CDF plots for DC1 and Tsth20	35
5.5	CDF plots for <i>T. harzianum</i> and <i>T. reesei</i>	36
5.6	CDF plots for <i>T. virens</i>	37
5.7	BUSCO counts	39
5.8	Example Secondary Metabolite Gene Clusters Identified by antiSMASH	43
5.9	Breakdown of identified regions	46
5.10	Example of a region with many gene calls	47
5.11	Predictions on opposing strands	48
5.12	Percentage of proteins with Pfam matches	49
5.13	Agreeing Pfam matches	51
5.14	Split Pfam matches	55
5.15	RefSeq absence with IPS evidence	56
5.16	Total tblastn hits	56
5.17	Results from tblastn by reference	58
5.18	Regions of agreement in AT-rich genomic sequence	58
5.19	Gene counts in AT-rich regions	59

List of Abbreviations

BGC	Biosynthetic gene cluster
BLAST	Basic local alignment search tool
BUSCO	Benchmarking universal single-copy orthologs
CDS	Coding sequence
CPU	Central processing unit
DNA	Deoxyribonucleic acid
GFF	General feature format
GC	Guanine cytosine
GMO	Genetically modified organism
HGT	Horizontal gene transfer
HMM	Hidden Markov model
Kb	Kilobases
Mb	Megabases
Mb	Mega-base
MITE	Minature inverted-repeat transposable element
NCBI	National center for biotechnology information
NGS	Next generation sequencing
PCR	Polymerase chain reaction
RNA	Ribonucleic acid
RSMI	Root, soil and microbial interactions
TE	Transposable elements
TDR	Terminal direct repeat
TIR	Terminal inverted repeat
UTR	Untranslated region
WGS	Whole genome shotgun (sequencing)

1 Introduction

Trichoderma species are fungi from the Hypocreaceae family that are commonly found in soil and aid in the decomposition of plant matter. They are also known for their ability to produce a variety of enzymes involved in a number of different roles, including cellulose degradation, nitrogen fixation, antibody production, and cellular signalling. Some of these species are also used in agriculture as biocontrol agents against plant pathogens. Some the enzymes responsible for these roles are produced in large quantities, making them of interest to researchers in both industry and academia. In particular, the elevated production of secondary metabolites in *Trichoderma* species has been shown to have a positive effect on plant growth and health, making them an attractive option for inoculation of agricultural crops. The genomes of *Trichoderma* species contain a wealth of information about the genes and proteins involved in these processes, but they must first be annotated to identify the genes and their functions, which is a non-trivial task. In addition, there are a number of different tools available for genome annotation, each with its own strengths and weaknesses, making it difficult for a researcher to choose the right tool for their needs. Gene finding tools are commonly compared to one another in common model organisms, but there is little work done comparing these tools in fungi, and even fewer in *Trichoderma* species. A focus on *Trichoderma* is important in this case, as the gene finders rely on gene models that may not be representative of the genes found in these species. The sequencing of two novel *Trichoderma* genomes, DC1 and Tsth20, provides an opportunity to compare the performance of several gene finding tools in these species. Results from this work may be of interest to researchers studying crop resistance and plant growth promotion, as well as those interested in the mechanisms at play behind other benefits of *Trichoderma* species. In addition, this work will provide insight into the performance of gene finding tools in fungi and the evolutionary processes at play in these species.

The main goal of this work is to compare the performance of several gene finding tools in the genomes of two novel *Trichoderma* species, DC1 and Tsth20. In addition to these two genomes, we also use the RefSeq assemblies and annotations of three other *Trichoderma* species, *T. reesei*, *T. harzianum* and *T. virens*, as an additional basis for comparison, with the goal of comparing these tools in a generic context, rather than in the context of a gold standard, as a gold standard is not available for these species. The exact research questions pursued in this work are explained in detail in Chapter 3.

The general structure of this thesis is as follows. In Chapter 2, we provide background information on *Trichoderma* species, genome assembly and annotation, and gene finding tools. As mentioned before, the research questions are explored in Chapter 3, while in Chapter 4, we describe the methods used to answer those questions. In Chapter 5, we present the results of our analyses, including comparisons of the assemblies

and annotations of the DC1 and Tsth20 genomes to those of *T. reesei*, *T. harzianum* and *T. virens*, as well as comparisons of the performance of the gene finding tools. Finally, in Chapter 6, we discuss the implications of our findings, as well as potential future work in this area. We make note that in this thesis, we refer to annotations from NCBI’s RefSeq database as ‘RefSeq’, although the annotations may have been generated using a variety of tools and methods.

2 Background

This chapter provides a brief background on *Trichoderma* and gene prediction tools, providing context for this research. Section 2.3 discusses the evolutionary mechanisms that may have led to the increased production of secondary metabolites in *Trichoderma* species, such as horizontal gene transfer, gene duplication, and transposable elements. Section 2.2 discusses the novel *Trichoderma* genomes that were sequenced and assembled as part of this work, and how they relate to the research questions. Section 2.4 discusses the structure of eukaryotic genes as well as a background of gene finding tools and methods, while Section 2.5 provides an overview of secondary metabolites and their importance in *Trichoderma* species. Finally, Section 2.6 provides an overview of comparative studies of gene prediction tools in *Trichoderma* species, and how this work aims to fill the gap in literature.

2.1 *Trichoderma* Species and their Features

Trichoderma species are a group of fungi in the *Hypocreaceae* family and are found in a wide variety of soils, being well-known for their use in biomanufacturing of cellulases and hemicellulases and also known for their roles as non-toxic, avirulent, opportunistic plant symbionts [1, 2]. Many of these features of *Trichoderma* can be attributed to the large number of secondary metabolites produced by these fungi, which are used to interact with plants and other organisms in the soil [3]. *Trichoderma* species differ in the number of secondary metabolites they produce, with some species producing more than others [3]. This variation in secondary metabolite production is of great interest to researchers, as it can provide insight into how these fungi can confer benefits to plants, and how they can be used in agriculture and industry. In addition, the study of secondary metabolites in *Trichoderma* species can help us better understand the evolutionary processes that have shaped *Trichoderma* species over time.

There are many mechanisms that organisms may leverage that result in gain or loss of gene function, and ultimately the evolution of a species. One well-studied mechanism is that of horizontal gene transfer (HGT), which is the process of acquiring genetic material from another organism, rather than through inheritance from a parent organism [4]. HGT is a common mechanism in bacteria, but it has also been observed in fungi, including *Trichoderma* species [4]. Another mechanism of interest is the duplication of genes, which can result in the production of more than one copy of a gene, leading to increased expression of that gene [4]. This is particularly relevant in the case of secondary metabolite production, as many *Trichoderma* species have been shown to have multiple copies of genes involved in secondary metabolite production [3]. Lastly, transposable

elements (TEs) are another mechanism of interest, as they can insert themselves into the genome and disrupt or enhance the expression of genes [4]. TEs have been shown to play a role in the evolution of *Trichoderma* species, particularly in the case of secondary metabolite production [4].

Interestingly, both HGT and TEs tend to present themselves in regions where GC content differs from the rest of the genome [4]. Abnormal GC content is also associated with other genomic features, such as centromeres [5] and repetitive regions [6], both of which may have an effect on the production of secondary metabolites. As a result, it is important to consider the GC content of a genome when studying *Trichoderma* species, which forms the basis of the Research Question 3.10. Some studies have shown that *Trichoderma* species contain very few, if any transposable elements, which is unusual for fungi [7]. It is hypothesized that TEs were lost in *Trichoderma* due to repeat induced point mutations (RIP), which is a process that occurs in fungi where TEs are silenced and eventually lost from the genome [7]. Genes responsible for RIP mutations also tend to target GC nucleotide pairs, converting them to AT nucleotides, which explains the low GC content observed in some regions of *Trichoderma* genomes [4]. These conversions usually occur when a C nucleotide is adjacent to a downstream A nucleotide.

Given the proximity of *Trichoderma* species to other fungi, bacteria, and plants in soils, it is possible that *Trichoderma* species acquired genes which are involved in response to antagonistic or sympathetic intercellular interactions from other organisms. While difficult to prove, studies have shown that HGT events have occurred in a number of fungal species [8]. It is also important to note that some *Trichoderma* species have higher numbers of genes involved in secondary metabolite production than others, indicating an evolutionary divergence, making comparative analysis of *Trichoderma* species and strains a useful tool for understanding the mechanisms behind secondary metabolite production [3].

2.2 Novel *Trichoderma* Genomes

Recently, two strains of *Trichoderma* have been identified in the prairie regions of Alberta and Saskatchewan. These two strains, named Tsth20 and DC1, have been found to have beneficial properties when used as an inoculant for plants in soils. Tsth20 has been classified into a clade with *Trichoderma harzianum*, and may act as a bioremediation tool in soils contaminated with hydrocarbon content. DC1, while not yet classified, is believed to confer salt and drought tolerance to plants in dry, salty soils. While bioremediation and resistance to drought tolerance have been investigated in other strains of *Trichoderma* [9], the study of new strains is still useful in the effort of discovering unique mechanisms used by *Trichoderma*. Little is known about these strains, so DC1 and Tsth20 were sequenced at the Global Institute for Food Security in an initial attempt to better understand the mechanisms at work in these genomes. While this research does not explicitly identify genomic elements related to the beneficial properties of their genomes, it may serve as a foundation for future research of *Trichoderma*. The assembly of these genomes come as a result of Research Question 3.1.

2.3 Evolutionary Landscape of *Trichoderma*

The evolutionary landscape of *Trichoderma* species is an important aspect of this work. DC1 and Tsth20 have not yet been classified into a specific *Trichoderma* species, although Tsth20 has been tentatively declared a strain of *T. harzianum*, and further understanding of their evolutionary relationships to other *Trichoderma* species may help explain how DC1 and Tsth20 originated. Figure 2.1 shows the phylogenetic relationships among selected *Trichoderma* species as well as the outgroup *Fusarium graminearum*. The tree was constructed using cursory alignments of the *Rpb2* protein, which is a commonly used protein for phylogenetic analysis in fungi [10]. Representative *Rpb2* proteins from the various fungal species were extracted from the NCBI resource for fungal orthologs, while representative *Rpb2* proteins from DC1 and Tsth20 were identified from the results generated later on in this work. *Trichoderma* species separate into a number of distinct clades, and these clades can vary between studies as evidenced by the low confidence in placement of *T. atroviride* and *T. gamsii*. Lack of confidence in evolutionary relationships is not uncommon in *Trichoderma* species, and is likely due to the high levels of horizontal gene transfer and gene duplication that have occurred in these species [4]. This makes *Trichoderma* species interesting to study but difficult to classify.

Phylogeny is also important when selecting genomes for comparative analysis of bioinformatics tools. Some tools may vary in their performance depending on the nature of the genome they are applied to, making a heterogeneous selection of genomes important for understanding the performance of the tools in different contexts. In this work, we selected individuals from three well-studied *Trichoderma* species, *T. reesei*, *T. harzianum*, and *T. virens*, which are all commonly used in research and industry, but are relatively distant from each other evolutionarily according to more detailed studies [10].

2.4 Gene Predictions and Complementing Similarity Searches

Gene prediction methods are generally based on the idea that genes can be identified by searching for patterns in a sequence, which match an expected gene structure or model. These gene structures begin from a 5' start codon, continue through a series of exons and introns, and end with a 3' stop codon [13]. Flanking the gene structure are promoter regions, which are typically found upstream of the start codon, and untranslated regions (UTRs), which are found both upstream of the 5' start codon and downstream of the 3' stop codon. Promoter regions are important for the regulation of gene expression, while UTRs are important for the stability and translation of the mRNA transcribed by the gene [13]. Sequences are translated from the 5' start codon through to the 3' stop codon, splicing together exons and ignoring introns. An example of a eukaryotic gene structure is shown in Figure 2.2.

Gene prediction methods vary in their approach, but they generally fall into two categories: *ab initio* methods and evidence-based methods [14]. *Ab initio* methods rely on the identification of patterns in the sequence that match an expected gene structure, while evidence-based methods use prior information such

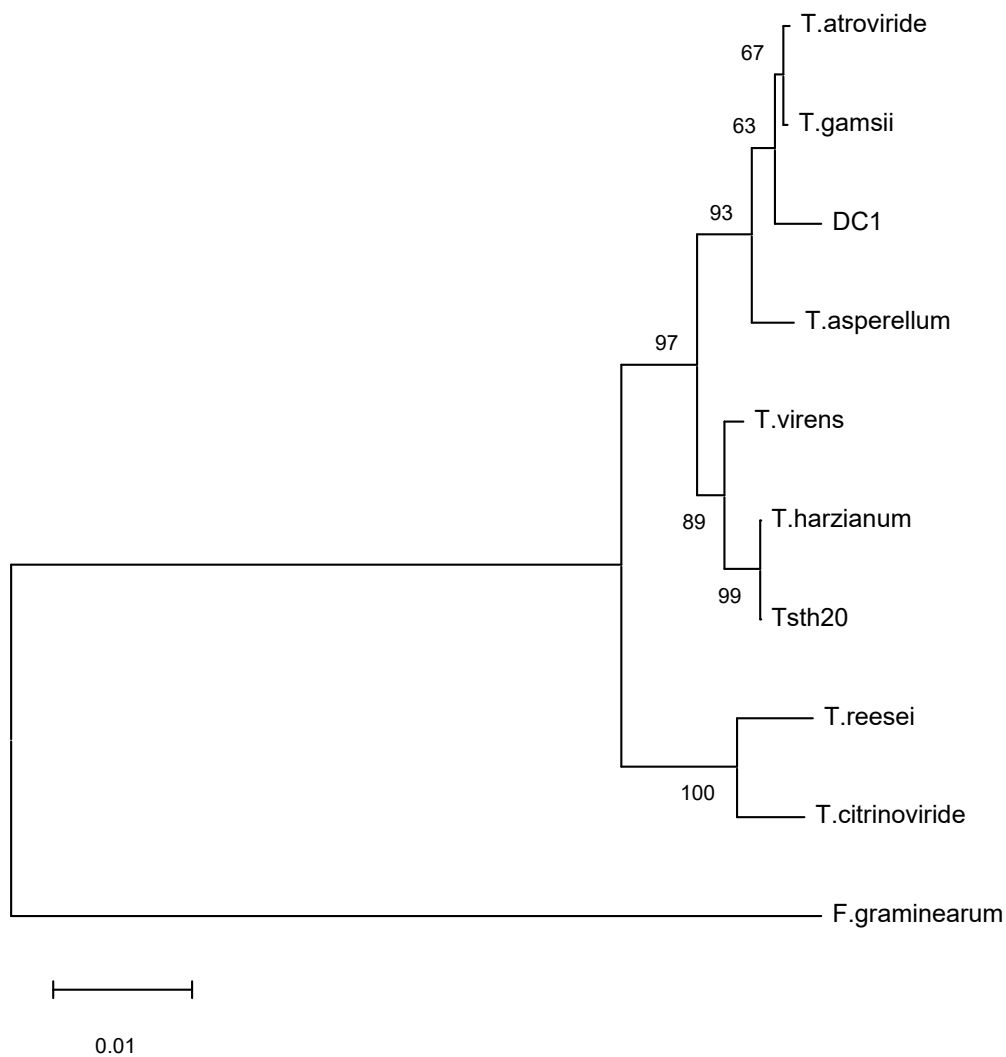


Figure 2.1: Phylogenetic tree illustrating relationships among selected *Trichoderma* species based on multiple sequence alignments of *Rpb2*. Evolutionary history was inferred using the Neighbour-Joining method [11]. The optimal tree with the sum of branch lengths = 0.146 is shown. The percentage of replicate trees in which the associated taxa clustered together in the bootstrap test (1000 replicates) are shown next to the branches [12]. The tree is drawn to scale, with branch lengths in the same units as those of the evolutionary distances used to infer the phylogenetic tree. An example branch of length 0.01 is shown.

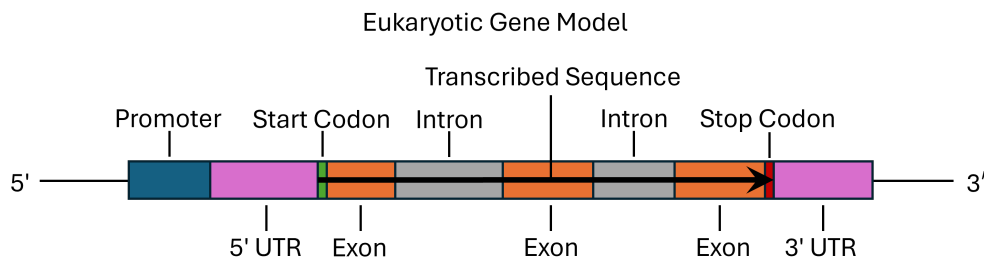


Figure 2.2: Example of a eukaryotic gene structure, showing a promoter region (blue), the 5' and 3' UTRs (purple), start codon (green), exons (orange), and introns (gray). The directed arrow indicates the direction of transcription.

as RNAseq data, expressed sequence tags (ESTs), and expressed protein sequences to identify genes within a new genome [14]. Gene prediction methods can also vary in their ability to identify different components of a gene, such as the promoter region, UTRs, and introns. Some methods may only identify the coding sequence (CDS) of a gene, while others may also identify the promoter region and UTRs [14]. Making things more complicated, different gene prediction tools may predict different combinations of exons and introns for a gene, and some tools may even predict multiple combinations of exons and introns, leading to multiple transcripts for the same gene. Given the already complex nature of secondary metabolites, it is important to consider the impact of variability in gene prediction methods, especially in the case of *Trichoderma* species and their secondary metabolite production. This consideration is one of the main motivators behind this research, and forms the basis of Research Question 3.9.

As mentioned earlier, gene finding tools generally fall into one of two categories: *ab initio* and evidence-based. *Ab initio* methods are often used when no prior information is available, and include tools such as GeneMark [15], GenScan [16], and GlimmerHMM [17]. Evidence-based methods are used when prior information is available [14], and include tools such as BLAST and other alignment tools to identify sequences in the genome that are similar to known sequences. More recently, hybrid tools have been developed that combine both *ab initio* methods and evidence from RNAseq data, ESTs, and protein sequences to identify genes in a genome [14]. These hybrid tools are often more accurate than either *ab initio* or evidence-based methods alone, as they can take advantage of the strengths of both approaches. Some examples of hybrid tools include AUGUSTUS [18], Braker2 [19], and MAKER [20]. GeneMark and Braker2 were selected as tools for comparison in this work as they are both widely used in the field, with high citation counts of over 2000 and 1600 respectively on individual papers, and also represent an *ab initio* method and a hybrid method, respectively. Braker2 was a particularly interesting option, as it combines GeneMark with AUGUSTUS as part of its workflow. Purely evidence-based gene prediction methods were not selected for this work, as they require prior information such as RNAseq data, which was not available for DC1 and Tsth20. The NCBI RefSeq annotations for the additional *Trichoderma* assemblies were also used as a point of comparison as they are generated using a combination of *ab initio* and evidence-based methods, and are considered to be high-quality annotations. RefSeq annotations may be generated by the NCBI Eukaryotic Genome Annotation

Pipeline [21] or by other sources as long as the annotation is approved by NCBI for the organism of interest. To simplify discussion surrounding the NCBI representative annotations, they will be referred to as RefSeq throughout the rest of this work.

While the basic structure of a gene may be present in a sequence, the gene might not code for a functional protein, or in the case of secondary metabolites, a component that helps produce them. To assess potential function of a gene product, gene prediction methods are often used in conjunction with similarity searches, which compare the predicted genes to known genes in other organisms to identify potential functions, and serve as a form of validation for the predicted genes [13]. Similarity searches can come in several different forms, such as BLAST searches, which compare the predicted genes to a database of known proteins, or InterProScan, which compares the predicted genes to a database of known protein domains and binding motifs [13]. These similarity searches can provide additional information about the predicted genes, such as potential functions, and can also help to validate the predicted genes by comparing them to known genes in other organisms. Similarity searches can also be used to evaluate the completeness of the predicted genes. One example would be comparing predicted genes to a set of conserved single-copy orthologs expected in fungal genomes, such as those provided by BUSCO [22]. Together, these methods can provide a more complete picture of the gene predictions and their potential functions, which allows us to better understand the performance of these gene finders in the context of fungal genomes. This forms the basis of Research Questions 3.6, 3.7, and 3.5.

2.5 Secondary Metabolites

While cellular products essential to an organism’s viability are of great interest, there are a vast number of cellular products that while not essential, still provide great benefit to the organism and may be necessary for survival in some situations [23, 3]. These other cellular products are known as secondary metabolites, and they are involved in several roles ranging from cellular signalling to antibiotic activity, making them a frequent subject of study in pharmaceutical research. Enzymes that produce secondary metabolites are comprised of non-ribosomal peptide synthetases (NRPS) and polyketide synthetases (PKS), which synthesize amino acids and other basic enzymatic building blocks into proteins [24]. Genes encoding the NRPSs and PKSs responsible for production of secondary metabolites are often found in clusters within a genome, but their products remain unknown [3].

The study of these products is difficult as the genes encoding the modules that make up NRPSs and PKSs are not expressed under normal laboratory conditions [3]. The NRPSs and PKSs that have been studied have been shown to be large enzymatic structures containing several functional modules, each responsible for a specific step in the synthesis of a protein [3]. Studying these complicated mechanisms and their genes requires as close to a complete set of gene predictions as possible, as the absence of one gene encoding a module could affect the entire protein complex. This concept of completeness of gene predictions is explored in Research

Question 3.5. In addition, gene predictions can be processed further with tools such as antiSMASH, which can identify secondary metabolite biosynthetic gene clusters (BGCs) in a genome [25]. Comparison of predicted genes with known secondary metabolite BGCs can provide insight into completeness of the predicted genes with a specific focus on secondary metabolite production. This work is not covered in this thesis, but is worth noting for context as it is a common approach in the field of secondary metabolite research.

Another interesting feature of NRPSs and PKSs is the length of the genes encoding them, which can be quite long, often exceeding 10,000 base pairs in length [24]. Examining the outputs of gene prediction tools can provide insight into the lengths of predicted genes, and whether the tools are able to capture the full length of these large genes. This topic serves as the basis for Research Question 3.4.

2.6 Gene Prediction Tools in Fungi

With the biological context laid out, we can now turn our attention to the computational aspect of this work. Many *Trichoderma* species have been assembled and annotated in the past, mostly for the purpose of understanding the mechanisms behind secondary metabolite production [3]. Gene predictions in this context are typically compared against sequences from other organisms to derive functional information, and to validate the predicted genes [13]. In contrast, very few studies have been conducted comparing gene prediction tools in fungal genomes, and even fewer have been conducted in *Trichoderma* species. This is surprising, given the large number of gene prediction tools available, and the fact that many of them are used in the field of fungal genomics [14]. Comparisons of gene prediction tools are common when new tools are released, but these comparisons are often limited to a few tools, and tend to be focused on achieving higher counts of interesting genes, rather than comparing the results in a biological context [26]. Complicating the matter further, the quality of genome sequences available may vary, making evaluation of different gene prediction tools difficult when working with multiple genomes and lesser-studied fungal species with low quality assemblies. Given the gap in literature regarding comparative analysis of gene prediction tools in *Trichoderma* species, and the availability of the novel DC1 and Tsth20 genomes, this work aims to fill that gap by comparing the outputs of several gene prediction tools on these two genomes along with three other frequently studied *Trichoderma* species.

Another pitfall of existing literature is that many studies compare tools only against a reference annotation, such as the RefSeq annotations for some organisms. However, these reference datasets are not always available, may not be representative of the genome being studied, and may be incomplete. In addition, it is helpful to compare the outputs of gene prediction tools against each other, as this can provide insight into the strengths and weaknesses of each tool relative to one another. For these reasons, this work aims to compare outputs of gene prediction tools in an all-to-all manner, rather than an all-to-one manner.

One positive commonality between studies involving gene prediction tools is the inclusion of more technical metrics of gene finders, such ease of installation and use, run times, memory usage, and other computational

metrics. These metrics can be important when selecting a tool, as users may be limited by the computational resources available to them, or may require results in a timely manner. Additionally, not all users are comfortable with command line tools, and may prefer a tool that is easy to use, and includes useful documentation.

3 Research Questions

With an ever-increasing number of gene prediction tools available, and the continuous discovery of new and unique genomic sequences, it is important to assess the performance of these gene finding tools. This chapter will address the following research questions, which will be answered in the context of *Trichoderma* genomes. The answers to these questions will provide insight into the performance of select gene finding tools, and may help inform users when selecting a gene prediction tool for use with fungal genomes and *Trichoderma* genomes in particular.

3.1 *Trichoderma* Assembly Results

Since gene finding tools operate on an assembled genomic sequence, it must follow that the results will be influenced by the supplied assembly. Before applying gene finders to the new assemblies, we should first evaluate the assemblies using general assembly metrics and compare with them existing assemblies. We ask: **how do assemblies of DC1 and Tsth20 compare to existing *Trichoderma* assemblies?** With these isolates being from the *Trichoderma* family, we expect assembly metrics to be similar in nature to existing assemblies from NCBI, but are they? In addition, what are the quantitative features of these assemblies?

3.2 Profiling of Gene Finding Tools

Different gene finding tools may predict different types of features associated with gene structures. The question arises: **which (if any) gene finders predict additional features outside the standard gene model?** Additional features in this case include promoter sequences, transcription binding sites, activating sequences and other upstream or downstream sequences. In addition, different gene finding tools employ differing programming languages and algorithms which raises several questions. **How are these gene finding tools implemented? Is the software straightforward to install? Are the tools user-friendly? What is the processing time and memory consumption of different gene finding tools in the context of *Trichoderma*?**

3.3 Number of Features Predicted

One common method for evaluating gene finding tools is by comparing the number of genes and coding sequences each tool predicts. This gives insight into the sensitivity of each tool, as well as a sense of

completeness of the predicted gene set. We will compare the number of features predicted by each tool, and ask: **do gene finders predict similar numbers of features in *Trichoderma* genomes?**

3.4 Lengths of Predicted Genes

Genome assemblies can contain a wide range of gene lengths, and for some users, genes of a specific length may be a key point of interest, as is the case in *Trichoderma*, where genes encoding non-ribosomal peptide synthases are generally quite large [24]. Thus, it is important to assess the lengths of predicted genes from different gene finding tools. We will compare the lengths of predicted genes from each tool and ask: **do gene finders predict genes of similar lengths in *Trichoderma*?**

3.5 BUSCO Completeness

The use of existing benchmarks for gene finding performance is useful when assessing performance of gene finding tools, particularly in the case of genes that should be evolutionarily conserved. **Do gene finders perform similarly when predicting conserved single-copy orthologs expected in fungal genomes?**

3.6 Validation of Predicted Genes via InterProScan

In an effort to validate, or at least provide supporting evidence for any given gene prediction we will apply InterProScan to coding sequences predicted by each of the gene finders to identify features associated with protein function. Genes will be considered as ‘valid’ if the gene’s protein sequence contains binding sites, motifs, or other functional characteristics of proteins. Using the results from InterProScan, we ask: **in the context of *Trichoderma*, to what extent do proteins predicted by gene finders contain functional signatures? Do gene finders predict similar numbers of genes with functional signatures?**

3.7 Similarity Searches of Predicted Genes with tblastn

In a similar manner to validation via InterProScan, we can use protein alignments to search for similarity between the protein products of predicted genes and existing protein sequences. This will allow us to assess the degree of similarity between proteins from predicted genes and known proteins, which can provide further validation of gene predictions and help infer possible function. We ask: **do gene finders predict genes whose translations are similar to known proteins in *Trichoderma*? Do gene finders identify comparable numbers of proteins that show similarity to known proteins?**

3.8 Identifying Candidate Secondary Metabolite Biosynthetic Gene Clusters

Secondary metabolites produced by *Trichoderma* species are of great interest due to their potential applications in agriculture and medicine. Identifying biosynthetic gene clusters responsible for the production of these secondary metabolites is crucial for understanding their biosynthesis and potential utility. We ask: **do gene finders differ in their ability to predict genes within secondary metabolite biosynthetic gene clusters in *Trichoderma* genomes?**

3.9 Agreement Between Gene Finders

With several sets of gene predictions, it is important to assess the agreement and disagreement between the predictions of different gene finding tools. To answer this question, we will identify ‘regions’ of overlapping predictions. A region is defined as a set of overlapping gene predictions, where a gene prediction from one tool overlaps with a gene prediction from another tool. With overlapping regions identified, we can determine any agreement or disagreement in predictions between gene finding tools from which we can ask: **do gene finders agree on their predictions? If no, to what extent do they disagree?**

3.10 Performance in AT-rich Genomic Sequence

Results from the assembly process show that AT-rich (GC-poor) sequence content is prevalent in most *Trichoderma* assemblies we selected. These regions of assemblies present an interesting opportunity to assess gene finding performance in regions of anomalous GC content, since features such as transposable elements and repetitive sequences are often found in these regions. We ask: **do gene finders perform differently in AT-rich regions of *Trichoderma* genomes? How well do gene finders perform in these regions?**

3.11 Selection of a Gene Finding Tool

With all the results generated, we can provide insight to the question: **when studying *Trichoderma*, which gene finding tool should one choose in any given context?**

4 Data and Methodology

4.1 Methodology Overview

The general processing steps performed in this work can be grouped into four processing stages, or workflows. A simplified overview of these stages and the flow of data between them is shown in Figure 4.1. In ascending order, each processing stage is discussed in more detail in Sections 4.2, 4.3, 4.4, and 4.5.1.

4.2 Sequence Processing and Assembly Workflow

The sequence processing and assembly workflow is shown in Figure 4.2. Processing steps are represented by boxes, results of interest generated by this work are represented by folders, and datasets that were not produced by this work are represented by cylinders. Other workflows that use datasets produced by this workflow are depicted as arrows outside the brown sub-graph. In this work, we note that both Illumina ©® and Nanopore ©® are copyrighted and registered trademarks of their respective owners. From here on out, we will refer to these technologies as Illumina and Nanopore without the associated symbols for brevity.

The sequencing data used in the assembling DC1 and Tsth20 were generated by Brendan Ashby as part of their MSc. work in the Biology department at the University of Saskatchewan under the study of Dr. Kaminskyj. Both Illumina and Nanopore sequencing technologies were used to generate genomic sequences for both organisms.

Illumina [27] sequences were first trimmed and filtered using Trimmomatic [28] to remove adapter content and low-quality sequences, as low-quality sequences may affect the assembly process, resulting in fragmented and erroneous assemblies. The processed Illumina sequences and Nanopore [29] sequences were then supplied to the hybrid genome assembly tool NextDenovo [30]. As Illumina sequences are of extremely high accuracy, especially after trimming, the processed Illumina data was then re-used to polish the assemblies with NextPolish [31]. Polishing is used because long-read sequences used in the assembly process generally have lower base-calling accuracy, at approximately 90.95%, than Illumina sequences at 99.9%, so the Illumina data can be used to correct any assembled sequences that may contain base-call errors. The resulting assemblies were then analyzed to identify AT-rich genomic sequence, analyzed with QUAST for general assembly metrics, and then used as input in the gene prediction workflow.

The processing steps Trimmomatic filtering, NextDenovo assembly, NextPolish Polishing and QUAST assembly assessment are described in Section 4.2.1, and were used to answer Research Question 3.1. The

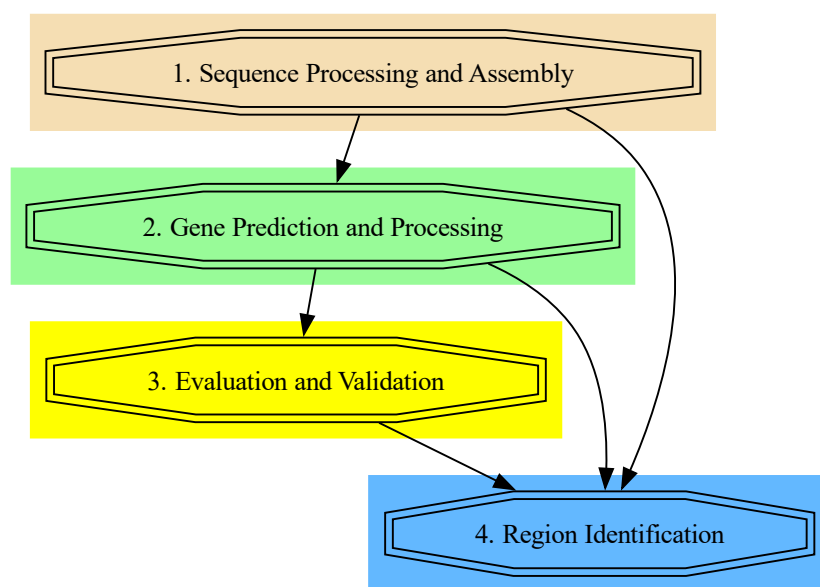


Figure 4.1: A simplified workflow used in this work. Each processing stage is represented by an octagon, and the flow of data are indicated by directed edges between them. In this figure and subsequent figures, the sequence processing stage is colored brown, the gene prediction and processing stage is colored green, the evaluation and validation stage is colored yellow, and the region identification process is colored blue.

AT-rich sequence identification process is described later in Section 4.2.2, and were used to answer Research Question 3.10.

4.2.1 Sequence Pre-processing and Assembly

Genomic sequences used for assembly were generated for DC1 and Tsth20 using Nanopore [29] and Illumina [27] sequencing technologies by Brendan Ashby at the Global Institute for Food Security at the University of Saskatchewan. Short-read sequences were first examined with the FastQC tool [32] (default parameters) to evaluate sequencing quality and general metrics associated with the sequencing run. The short-read sequences were then trimmed to remove adapter content or low-quality sequences, as these may lead to fragmented assemblies and erroneous base calls. Illumina sequencing data was filtered using Trimmomatic v0.38 [28] with the following criteria: SLIDINGWINDOW:4:28, LEADING:28, TRAILING:28, MINLEN:75. Only surviving paired-end reads were used for further analysis. After filtering, 5.57 million paired-end Illumina reads between the length of 75bp and 250bp remained for DC1, while Tsth20 had 7.94 million paired-end Illumina reads between the length of 75bp and 250bp. Single reads that lost their pair during filtering were discarded.

Nanopore sequencing for DC1 and Tsth20 was also performed by Brendan Ashby. For DC1, 6.25 million reads were obtained with a mean quality of 9.9 and a mean length of 1,062bp. For Tsth20, 1.27 million reads were obtained with a mean quality of 10.2 and mean length of 4,788bp. All Nanopore reads used were of Nanopore’s Q quality score of 7 or higher, and no additional filtering or trimming was performed.

It has been shown that assemblies that use both short-read and long-read sequencing data during the assembly are of higher quality and less fragmented than assemblies using only one form of sequencing. Genomic Nanopore and Illumina sequences from DC1 and Tsth20 were assembled into contigs using the hybrid assembly tool NextDenovo [30] v2.5.0 with default parameters. Although hybrid genome assemblies are of higher quality, they may still contain base calling errors after assembly. The highly accurate short-reads can be used to correct base calling errors after assembly. To do this, the hybrid assemblies were polished with Illumina sequences using NextPolish [31] v1.4.1 and default parameters. Final assembly metrics were calculated using QUAST v5.0.2 [33] with default parameters.

4.2.2 Identifying AT-rich Genomic Sequence

Initial exploration of the Illumina sequence data generated for DC1 and Tsth20 showed bi-modal distributions of GC content in the reads. An example of the distribution of GC content from Illumina sequences for DC1 is shown in Figure 4.3. We can see that there is an increase in AT-rich (GC-poor) sequences around 16% GC content, which is unexpected in comparison to the expected peak around 50% GC content. Because of this observation, we extended the GC content analysis further to identify regions of AT-rich genomic sequence in both DC1, Tsth20, and the other assemblies discussed in Section 4.3.1. This process is described below.

AT-rich genomic sequences were identified using sliding windows of GC content over each *Trichoderma*

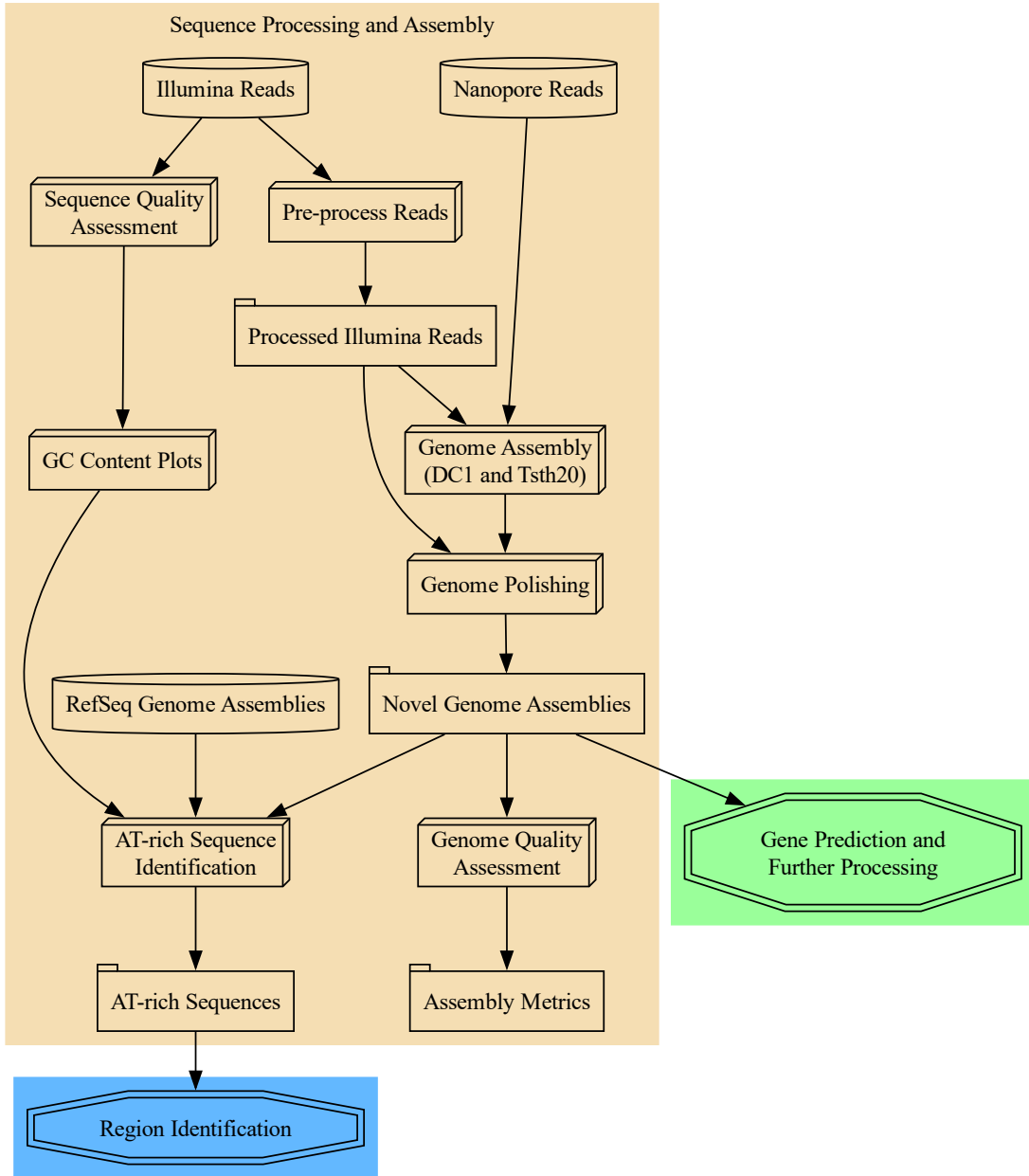


Figure 4.2: A workflow (in brown) depicting the steps taken to process input sequences, generate assemblies, and post-process those assemblies for use in other workflows. Processing steps are represented by boxes, datasets of interest that were produced by this work are represented by folders, and datasets that were not produced in this work are represented by cylinders. Other workflows are represented by green and blue arrows outside the brown graph. Directed edges indicate the flow of data.

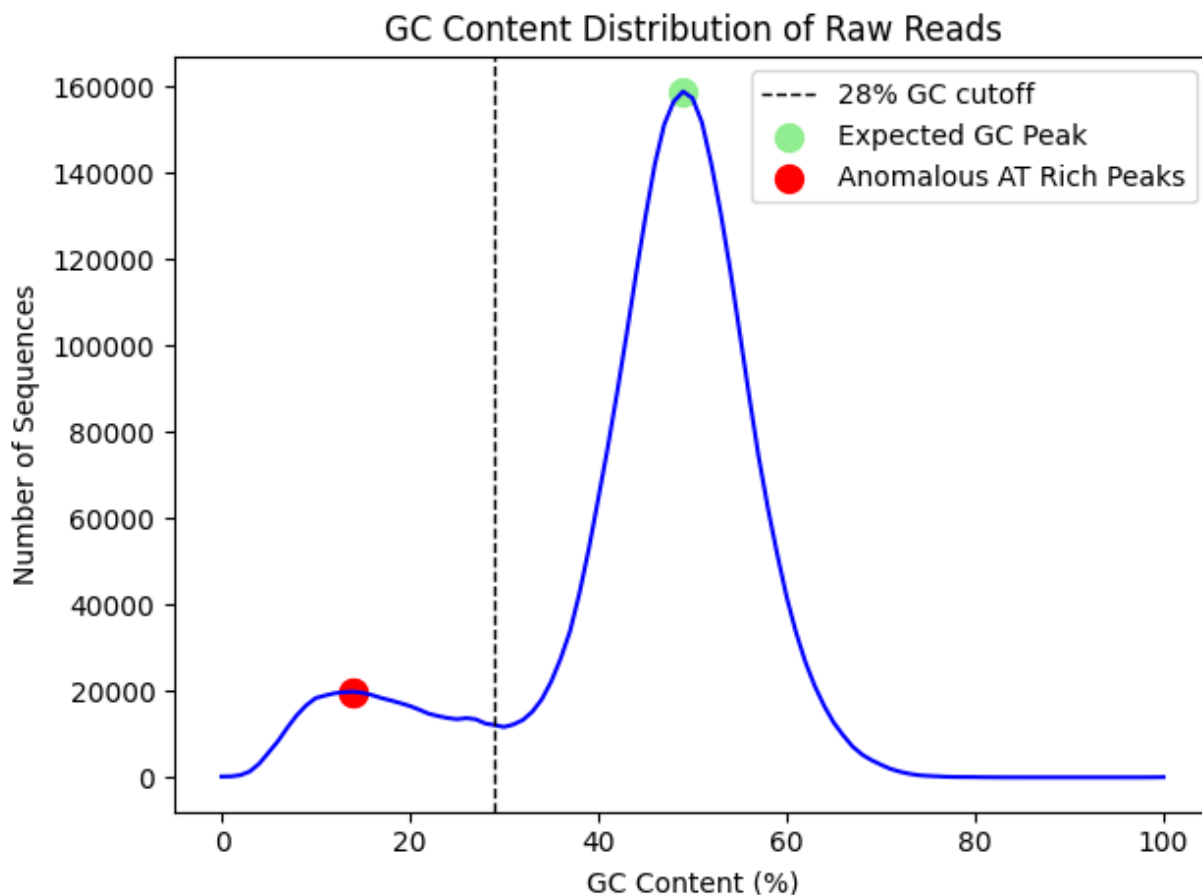


Figure 4.3: A plot of GC content of Illumina sequences generated for DC1. The X-axis indicates mean GC content of sequences while the Y-axis indicates the total number of sequences for a given mean GC content. The green circle highlights an expected peak around 50% GC content, while the red circle highlights an anomalous peak around 16% GC content. The dashed black line indicates the 28% GC content cutoff used to identify AT-rich sequences.

assembly. Sliding windows were generated using the isochore program from the EMBOSS tool suite v6.6.0 [34]. The sliding window used was 250bp wide with a 50bp shift. A window was classified as AT-rich if the sequence contained lower than 28% GC content, which was based on the observed distribution of GC content in Figure 4.3 and the overall trends observed in the other assemblies. The genomic sequences from AT-rich windows were mapped to a GFF file and processed using the region identification algorithm described in Section 4.5.1.

4.3 Overview of Gene Prediction and Further Processing

To better understand the methods used in gene prediction and further processing, the steps have been assembled into a pipeline shown in Figure 4.4. Processing steps are represented by boxes, results of interest generated by this work are represented by folders, and datasets that were not produced in this work are

represented by cylinders. Other workflows that use datasets produced by this workflow are depicted as arrows outside the green sub-graph.

First, RefSeq assemblies and their associated annotations of closely related *Trichoderma* species were selected from the National Center for Biotechnology Information (NCBI) for training and comparison. Both *ab initio* (GeneMark [15]) and evidence-based (Braker2 [19]) gene finders were selected for comparison to evaluate gene prediction behavior of both methods in the context of *Trichoderma* genomes. Gene prediction was performed using both gene finders on the RefSeq genomes and the genomes of DC1 and Tsth20. Coding sequences output by the gene finders were then examined to better understand distributions of CDS lengths produced by different gene finders. Finally, all gene predictions were then supplied to the evaluation and validation stage as well as the region identification stage.

Datasets from the National Center for Biotechnology Information (NCBI) were selected for training and comparison, and are further discussed in Subsection 4.3.1. The GeneMark gene prediction process is described in Subsection 4.3.2, while the Braker2 gene prediction process is described in Subsection 4.3.3. The results from these first two processes were used to answer Research Question 3.3. While performing gene prediction with Braker2 and GeneMark, observations on installation, usability, and performance were also recorded and used to answer Research Question 3.2. The process used to examine CDSs is described in Subsection 4.3.4 and the results were used to answer Research Question 3.4.

4.3.1 Selection of Datasets from NCBI

When evaluating gene-finders, it is important compare gene predictions to a standard, or control dataset, as a method of ground-truthing. Such datasets include reference genome sequences and their associated gene predictions from the NCBI as used in this work. In addition to comparing against control datasets, RNAseq data can also be used as training data for various tools, such as the evidence-based training used by Braker2. In this work, four forms of these inputs were selected from NCBI for comparison and training. The first two datasets selected were assemblies and associated gene predictions of three RefSeq *Trichoderma* organisms, with those organisms being *T. virens*, *T. harzianum*, and *T. reesei*. NCBI RefSeq accession IDs for those assemblies are GCF_000167675.1_v2.0, GCF_003025095.1_Triha_v1.0, and GCF_000170995.1_TRIVL_v2.0, respectively. *T. reesei* was selected as it is the most well-studied *Trichoderma* species, and widely considered to be the gold-standard for *Trichoderma* genomes [35]. *T. harzianum* and *T. virens* were also selected as additional datasets due to their prevalence in literature. We note that these genomes were assembled before 2018 using Illumina and Sanger sequencing data. No long-read sequencing was used in generating these assemblies. It is important to note that while RefSeq is considered a high-quality dataset, we did not use it as a ground truth in this case, but rather an additional point of comparison as it is possible that RefSeq annotations may also contain errors or omissions. This aligns with the method of region identification used later in this work, which is aimed to be a generic and adaptable method of comparing gene predictions, or any feature in GFF3 format, from multiple sources.

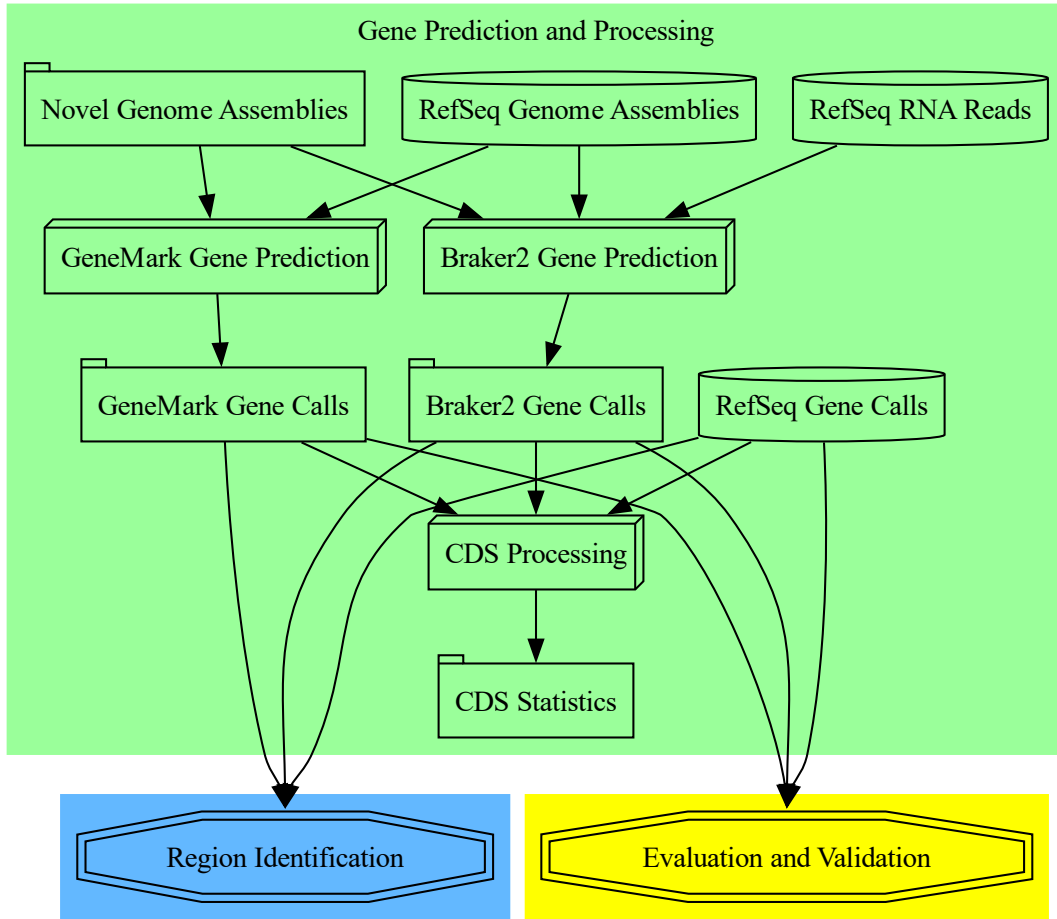


Figure 4.4: A workflow (green) depicting the steps taken to predict genes for use in other workflows and process CDS sequence lengths. Processing steps are represented by boxes, datasets of interest that were produced by this work are represented by folders, and datasets that were not produced in this work are represented by cylinders. Other workflows are represented by blue and yellow arrows outside the green graph. Directed edges indicate the flow of data.

The third set of selected inputs were RNAseq datasets from *T. reesei*, used as input in the training process for Braker2. The use of training data in gene finding applications is to produce a predictive model that factors in experimental evidence for a ‘tailored’ gene prediction model. These models can predict genes with more confidence and predict the correct gene structure more often [36]. The SRA accession IDs for sequences used are: SRR5229930, SRR8329344, SRR8329345, SRR8329346, and SRR8329347. Limited RNAseq datasets were used as there were few RNAseq datasets available at the beginning of this project.

The final set of inputs were protein sequences from RefSeq for the species *T. atroviride*, *Fusarium graminearum*, and *Saccharomyces cerevisiae*, which were for downstream validation of gene predictions. The NCBI accession IDs for these assemblies and associated protein sequences are GCF_000171015.1, GCF_000240135.3, and GCF_000146045.2, respectively.

4.3.2 Gene Prediction Using GeneMark-ES

Ab initio gene finding was performed using GeneMark-ES v4.71 [15] with DC1, Tsth20 and selected RefSeq *Trichoderma* assemblies used as input. Default parameters were used in all cases except for the inclusion of the fungal option, which allowed the use of a GeneMark model specific to fungal genomes.

4.3.3 Gene Prediction Using Braker2

For evidence-based gene finding, Braker2 v3.0.2 [19] was used with *Trichoderma reesei* selected as the reference organism for training. Illumina [27] RNAseq datasets from *T. reesei* were downloaded from the NCBI short-read archive using the SRA-toolkit [37] and were not trimmed or filtered prior to their use in Braker2 training. Default parameters were used to train a Braker2 model on *T. reesei* data except the addition of the fungal option, which was used for this analysis for improved gene-finding performance in fungi. The trained Braker2 *T. reesei* prediction model was then applied to all assemblies with default parameters.

Braker2 also provides a UTR option to predict upstream sequence features, but that option is experimental and was left off for this work.

4.3.4 Statistical Analysis of CDS Lengths

The distribution of gene lengths predicted by a gene-finder is an important feature to consider when selecting a gene finding tool, as genes of a particular length may be of interest to researchers, but more importantly, gene-finders may be biased to certain length distributions based on their underlying models. These biased distributions may differ from the true distribution of CDS lengths leading to propagation of suboptimal gene-finding results and results that may not reflect the true distribution of gene lengths. To determine if differences exist between distributions of predicted gene lengths from each tool, the \log_{10} lengths of coding sequences (CDS) were subject to a two-sample Kolmogorov-Smirnov test [38]. The null hypothesis of this test is that the variation between distributions of frequencies of gene lengths between two gene finders

can be explained by random chance, while rejecting the null hypothesis means that the samples originated from different distributions. This test does not explain the nature of the differences, only that a difference exists, and serves as evidence for further analysis.

4.4 Overview of Evaluation and Validation of Gene Predictions

Following the gene prediction process, prediction results can then be evaluated and used as queries against existing datasets as a form of validation. An overview of the evaluation and validation workflow is shown in Figure 4.5.

Coding sequences from the previous stage were evaluated and validated using several tools. Firstly, protein products from predicted CDSs were processed with InterProScan to identify Pfam domains as evidence for the validity of gene predictions. Secondly, CDSs were analyzed with BUSCO to assess prediction of conserved fungal genes. Thirdly, RefSeq proteins from related fungal species were used in tblastn searches of the five *Trichoderma* assemblies selected for the main body of this work. Finally, gene predictions from Braker2 and GeneMark were supplied to the antiSMASH tool for the identification of secondary metabolite biosynthetic gene clusters.

BUSCO evaluation, InterProScan Annotation, and BLAST alignments processes are described in more detail in Subsections 4.4.1, 4.4.2, and 4.4.3, respectively. The mentioned subsections align with Research Questions 3.5, 3.6, and 3.7, respectively. Identification of secondary metabolite biosynthetic gene clusters using antiSMASH is described in Section 4.4.4 and aligns with Research Question 3.8. Results from some of these processes are also used as inputs in the region identification workflow later.

4.4.1 Evaluating Gene-Finder Performance Using BUSCO

While novel gene predictions are of great interest, it is also important to confirm that gene-finders are predicting genes that are expected to be present in the organism of interest [39]. To determine if this is true, BUSCO v6.0.0 [40] was selected as a tool to determine a gene-finder’s ability to predict a set of conserved orthologous genes. The BUSCO Metaeuk pipeline was run using default parameters with the BUSCO sordariomycetes.Odb12 fungal lineage dataset to capture highly conserved genes in fungal genomes closely related to *Trichoderma* species. This pipeline was applied to all genome assemblies used in this work.

4.4.2 Functional Annotation Using InterProScan

The presence of a functional component in a predicted gene’s protein product is good evidence in favor of a gene-finder’s performance. To identify functional components, the protein products from each gene finder’s predictions were functionally annotated using InterProScan v5.6597.0 [41] without the match lookup service and with the following analyses included: AntiFam7.0, CDD-3.20, Coils-2.2.1, FunFam-4.3.0, Gene3D-4.3.0, Hamap-2023.01, MobiDBLite-2.0, NCBIfam-13.0, PANTHER-18.0, Pfam-36.0, PIRSF-3.10, PIRSR-2023_05,

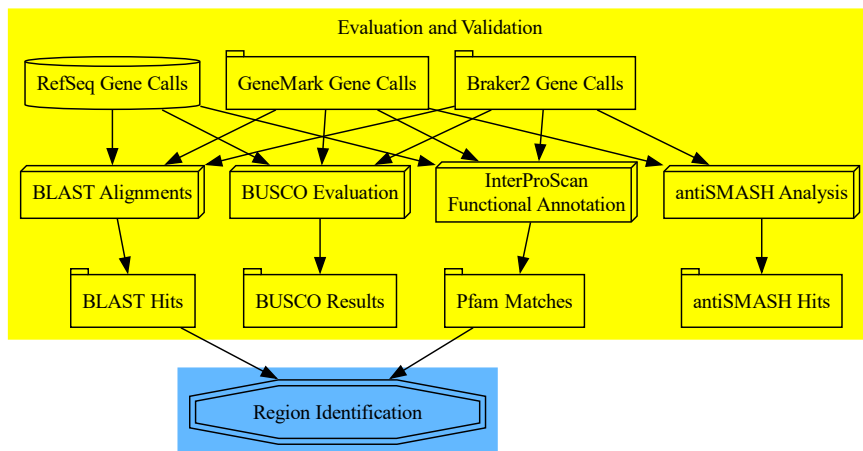


Figure 4.5: A workflow (yellow) depicting the steps taken to evaluate and validate predicted genes. Processing steps are represented by boxes, datasets of interest that were produced by this work are represented by folders, and datasets that were not produced in this work are represented by cylinders. Other workflows using results from these processing steps are represented by a blue arrow outside the yellow graph. Directed edges indicate the flow of data.

PRINTS-42.0, ProSitePatterns-2022_05, ProSiteProfiles-2022_05, SFLD-4, SMART-9.0, SUPERFAMILY-1.75. InterProScan was run on all predicted proteins using default parameters. The resulting protein annotations were mapped back to their genome assemblies and processed using the region identification algorithm in Section 4.5.1.

4.4.3 Ground-truthing with the Basic Local Alignment Search Tool (BLAST)

The *Trichoderma* assemblies chosen for comparison were used as the subject sequences in tblastn [42] searches with the query proteins discussed in Section 4.3.1. Default parameters were used in the tblastn searches. Hits from the tblastn alignments were then filtered for alignments with a minimum of 30% query coverage and 30% identity. These criteria were selected as they filter out spurious alignments while also retaining short alignments that capture functional or conserved protein coding sequence. Remaining alignments were then processed using the region identification algorithm in 4.5.1.

4.4.4 Identification of Secondary Metabolite Biosynthetic Gene Clusters Using antiSMASH

Secondary metabolite biosynthetic gene clusters (BGCs) are groups of co-located genes that work together to produce secondary metabolites, which are compounds not directly involved in normal growth, development, or reproduction of an organism. Identification of BGCs is important for understanding the metabolic potential

of an organism and for discovering novel natural products. Due to time constraints, the following analysis was only performed on the Tsth20 assembly, as it contains the fewest number of contigs and has been observed interacting with plants in soil. To identify candidate BGCs in the Tsth20 assembly, the antiSMASH v7.0.0 fungal version [25] was used with default parameters, using the Tsth20 genome assembly and associated gene predictions as inputs. The outputs of antiSMASH are HTML reports and GenBank files containing identified BGCs and were used in cursory analysis of potential secondary metabolite production in Tsth20. The use of default parameters omits the use of additional similarity methods, which results in some empty fields in the output reports. However, this analysis was only intended to identify confirm the presence of BGCs in Tsth20 and not to perform an in-depth analysis of their potential products.

4.5 Region Identification Overview

A major portion of this work was the process of identifying regions of overlapping features as a method of comparing gene calls from different gene prediction tools, which is a critical step in answering Research Question 3.9. The definition of a region and the algorithm used to identify them is described in detail in Section 4.5.1, and an overview of the region identification workflow is shown in Figure 4.6. The top of the graph shows the gene calls produced by the workflow described in Subsection 4.3, which are used as the main input to the region identification process. The inputs and resulting regions of other previously generated datasets is represented within the rounded rectangles. These additional datasets may contain other information such as BLAST alignments or functional annotations and can be included as features in the region identification process. Regions output by the region identification process were then visualized with the Integrative Genomics Viewer, described in Section 4.5.3.

4.5.1 Region Identification Algorithm

The gene predictions from Subsections 4.3.1, 4.3.2, 4.3.3 and the additional datasets from Subsections 4.2.2, 4.4.3, 4.4.2 were processed into regions of overlapping gene predictions. A region is defined as a pair of genomic coordinates containing one 5' coordinate and a subsequent 3' coordinate, between which are one or more overlapping features. We define a feature as a pair of genomic coordinates containing one start and one stop coordinate that is associated with a gene prediction, annotation, or other information of interest.

Overlapping features were processed into regions using a concept similar to Allen's interval algebra [43]. The algorithm first sorts all features on each contig by start coordinate and then iterates over each feature, identifying overlaps with previous features to identify regions of overlapping features. We approached this as a graph problem, where features make up a set of nodes V , and spatial overlaps between features make up a set of edges E , forming the graph $G = \{V, E\}$. This algorithm performs a transitive closure R^+ of G , resulting in sub-graphs, or regions of features in which each node has a path to each other node in that sub-graph. Pseudocode for the algorithm is shown in Algorithm 1. This processing was performed with

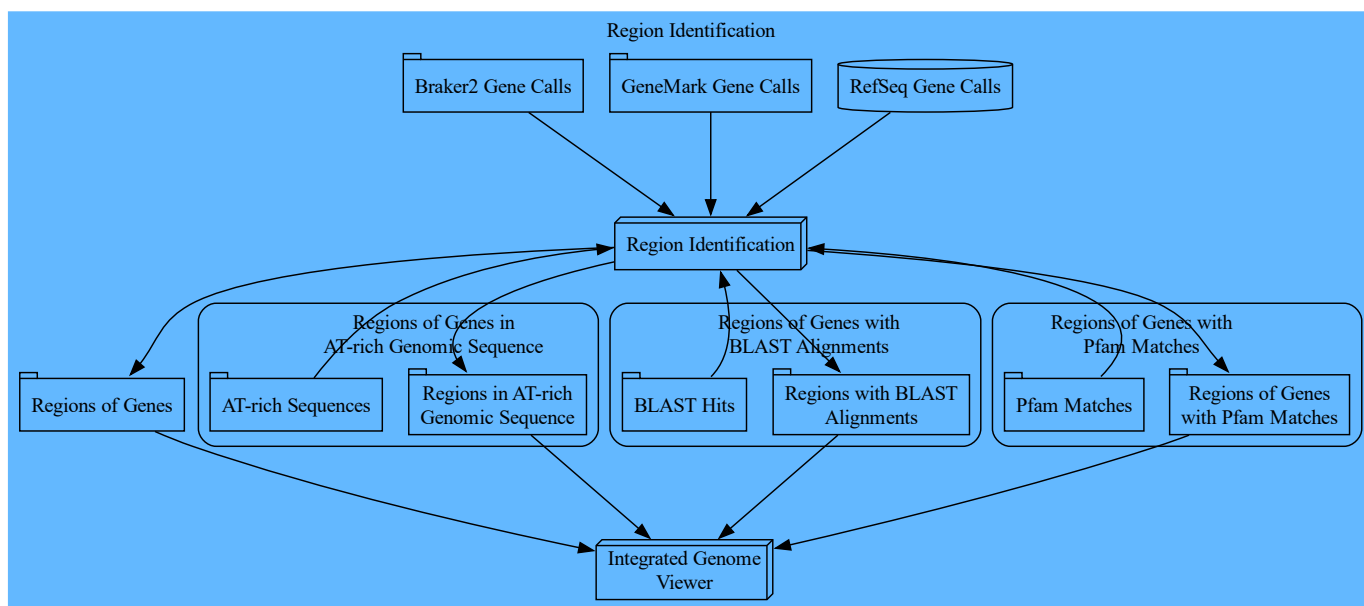


Figure 4.6: An overview of the region identification process for different datasets. Sections of the graph are outlined with rounded rectangles, indicating that the datasets within were processed along with gene calls and discussed in their own results sections. Processing steps are represented by boxes, datasets of interest that were produced by this work are represented by folders, and datasets that were not produced by this work are represented by cylinders. Directed edges indicate the flow of data.

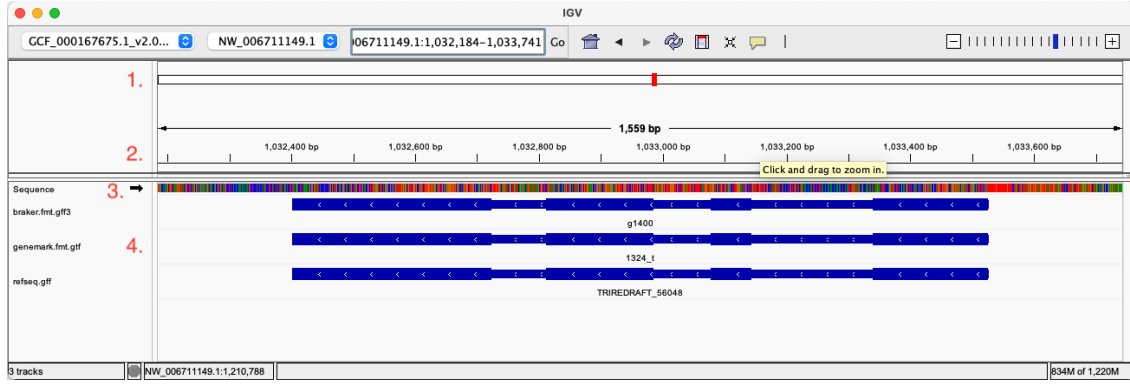


Figure 4.7: An example of an IGV screenshot used in portions of the results from this work. The top half of the image displays information about the reference sequence with track one showing the position of the viewer relative to the reference sequence, and track two showing the numerical positions of the reference. The lower half of the image displays more tracks in the form of additional layers of information mapped to the reference sequence. Track three displays the nucleotide at each position if the scale of the viewer is small enough. Track four comprises several subtracks, each containing features, the regions in blue, from a GFF file. The tracks in this example are limited to gene predictions, although tracks in other IGV screenshots may include features from other tools such as InterProScan and are discussed in their respective sections. Individual tracks are referred to in increasing numerical order in the track list beginning after the nucleotide track.

Python v3.12.4 [44] and the GFF package from BCBio [45] for easy GFF parsing. The resulting regions were then further processed to identify trends and behaviors of gene finders in the context of other annotated information and features.

4.5.2 Regions in AT-rich Genomic Sequence

AT-rich genomic sequence regions were identified using the process described in Subsection 4.2.2 and were used input to the region identification process.

4.5.3 Visualization of Identified Regions

Results from the region identification process were visualized using the IGV genome browser v2.19.1 [46]. An example of an IGV screenshot is shown in Figure 4.7.

Algorithm 1 The general algorithm underlying the region identification process.

```
INPUT: list of GFF files, GFFList
RETURNS: list of regions, regionList

for <each GFF in GFFList> do
    allFeatures  $\leftarrow$  GFF.features
end for

sortedFeatures  $\leftarrow$  sortStartCoord(allFeatures)
regionList  $\leftarrow$  []
firstFeature  $\leftarrow$  firstItem(sortedFeatures)
lastFeature  $\leftarrow$  lastItem(sortedFeatures)
for <each feature in sortedFeatures> do
    if feature = firstFeature then
        currentRegion  $\leftarrow$  newRegion(feature)
        currentRegion.start  $\leftarrow$  feature.start
        currentRegion.end  $\leftarrow$  feature.end
    else if feature = lastFeature then
        regionList.append(currentRegion)
        return regionList
    else if feature overlaps currentRegion then
        currentRegion.updatePositions(currentRegion, feature)
        currentRegion.addFeature(feature)
    else
        regionList.append(currentRegion)
        currentRegion  $\leftarrow$  newRegion(feature)
        currentRegion  $\leftarrow$  feature.start
        currentRegion  $\leftarrow$  feature.end
    end if
end for
return regionList
```

5 Results and Discussion

5.1 Overview

Results from this work are ordered according to the sequence in which data was processed, following the flow of Figures ?? . Since there are novel *Trichoderma* assemblies included in this work, the results begin with an overview of the assemblies of DC1 and Tsth20 in Section 5.2. Following this is Section 5.3, which discusses the installation, application and other features associated with the gene finders selected for this work. Results of the initial gene finding process are then presented in Section 5.4 followed by more detailed analysis of predicted CDS lengths in Section 5.5. Gene predictions from each gene finder and assembly combination are supplied to the BUSCO metric process in Section 5.6. Gene predictions from Braker2 and GeneMark in Tsth20 were also supplied to antiSMASH to identify candidate secondary metabolite biosynthetic gene clusters in Section 5.7. Results from the initial spatial comparison of gene predictions using the region identification process are then presented in Section 5.8. The region identification process is then applied with additional information from tblastn alignments, InterProScan annotation, and GC sequence composition in Sections 5.9, 5.10 and 5.11, respectively. A final comparison of all results and selection of an appropriate gene finder are then discussed in the Conclusion 6.1. The research questions behind these analyses are described in detail in Chapter 3.

RefSeq assemblies in these results are referred to by their scientific name for context and not by their associated NCBI accession IDs. Instead, NCBI accessions for the RefSeq assemblies and their associated annotations are listed here: *T. reesei* - GCF_000167675.1.v2.0, *T. harzianum* - GCF_003025095.1.Triha.v1.0, *T. virens* - GCF_000170995.1.TRIVI.v2.0.

5.2 Assemblies of DC1 and Tsth20

For general assembly metrics of DC1 and Tsth20, the QUAST [47] tool was used. Results from QUAST are shown in Table 5.1, from which we can make several observations. In DC1 and Tsth20, the total contig counts are an order of magnitude smaller than the NCBI RefSeq assemblies, indicating highly contiguous assemblies from nextDenovo [30] and nextPolish [30]. This is likely due to the use of Nanopore sequencing in the assemblies of DC1 and Tsth20. The total assembled lengths of DC1 and Tsth20 are similar when compared to RefSeq assemblies, ranging from 38Mb to 42Mb, except in the case of *T. reesei*, which is known to have a significantly smaller genome length [2], at roughly 33Mb.

Strain	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

Table 5.1: General assembly metrics produced by QUAST[47] (a genome quality assessment tool).

The largest contig size for each assembly varies greatly. DC1 and Tsth20 have the largest contigs of all assemblies being considered, which is again likely due to the inclusion of long-read sequencing data in the assembly process. The N50 values for all assemblies are above 1Mb, with DC1 and Tsth20 N50s being at minimum two times larger than other assemblies. N50 lengths nearing chromosomal lengths suggests that assemblies of DC1 and Tsth20 are better assembled than the other *Trichoderma* assemblies. While chromosome-scale contigs are not the sole indicator of genome quality, it does provide confidence in the quality of the input data and resulting assemblies. In general, the assemblies of DC1 and Tsth20 are of similar length to existing *Trichoderma* assemblies and the number of contigs reported match the number of ‘chromosome’ scale contigs reported in other work [2].

During initial investigation of the inputs to the assembly process, we observed that the Illumina reads have a bimodal distribution of GC content as shown in Figure 4.3. To see if this observation extended to the assemblies as well, 250 bp sliding windows were used to calculate GC content for all five assemblies included in this analysis. The results of this analysis are shown in Figure 5.1. AT-rich sequences are visualized by the left peak (windows containing $\geq 72\%$ AT nucleotides) of the distributions. Of the included assemblies, AT-rich sequences were identified in DC1, Tsth20, *T. reesei* and *T. harzianum*, with *T. virens* deviating from the other assemblies showing very few AT-rich windows. This lack of AT-rich sequence in *T. virens* nucleotide composition may indicate potential assembly issues or even misclassification of the organism. In addition to the confirmation of increased AT-rich sequence content in most assemblies, it appears that the distribution of GC content in *T. reesei* differs from the other assemblies. The curve of GC content for *T. reesei*, visualized in green in Figure 5.1, is shifted to the right of the other *Trichoderma* assemblies, indicating a higher overall GC content. We can not explain exactly why *T. reesei*’s curve differs so much at this time, but it is likely due to the much smaller *T. reesei* assembly and its properties. Investigation of these AT-rich sequences is continued in Section 5.11, where only sequences containing greater than 72% AT nucleotide content are considered, as the distributions begin to deviate from expected.

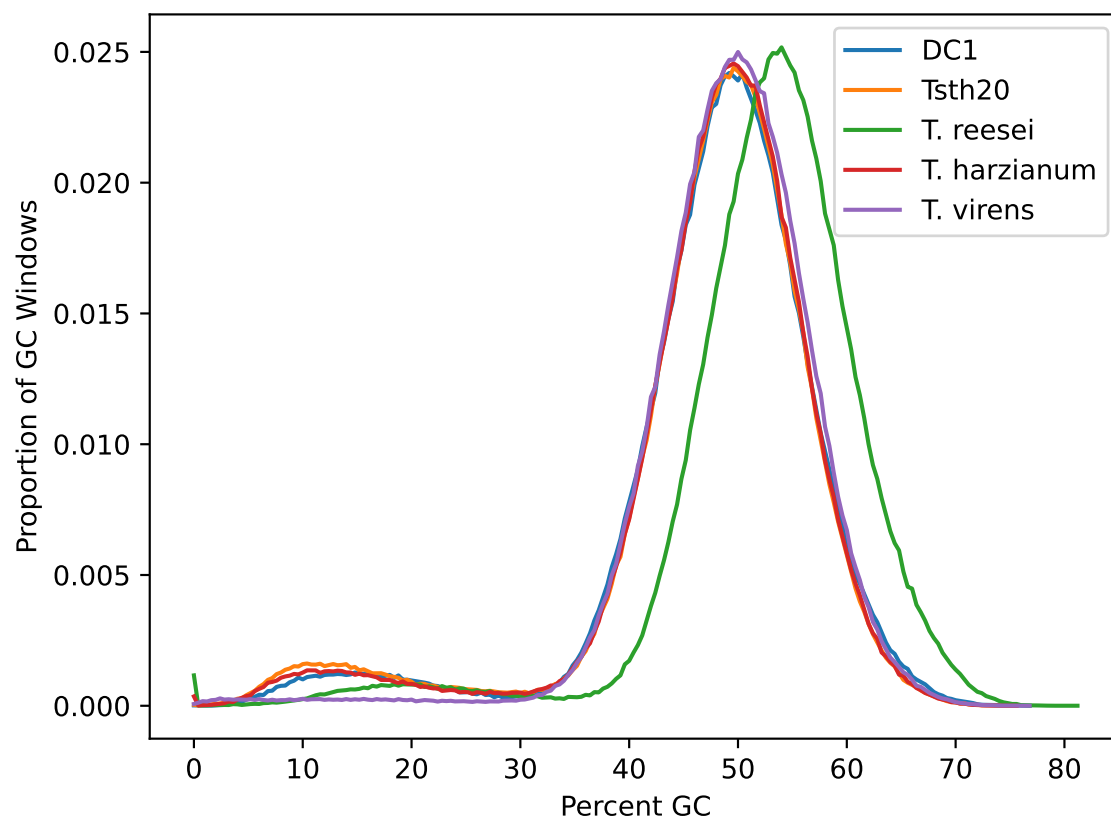


Figure 5.1: Plots showing proportions of sliding windows of GC content for each *Trichoderma* genome assembly.

5.3 Installation and Profiling Gene Finding Tools

While Braker2, GeneMark and RefSeq all predict coding sequences of possible genes for a provided input, the implementation of each tool is different, requiring more or less effort to install and run than other tools. This section discusses the implementations, installation procedures and execution of GeneMark and Braker2 in more detail. For more information on the versions and parameters used in this processing, please see Section 4.3.2 and Section 4.3.3. The computing platform used in this research is hosted and managed by University of Saskatchewan services, which is modelled around the HPC platform used by the Digital Research Alliance of Canada and software is managed similarly.

GeneMark [15] is a gene finding tool developed by the Georgia Institute of Technology with installation packages prepared for Linux and macOS. It is provided as licensed product in the form of a package which can be downloaded from their website after submitting a form. Once the necessary information is submitted, the user is provided with a key that must be placed in the appropriate location once the software is downloaded and unpacked. The core controlling methods of GeneMark are written in Perl, accompanied by several Python scripts and compiled executables. GeneMark was tested by the developers with Perl version 5.10 and Python 3.3. A number of Perl dependencies are also required, which can be installed via CPAN. The user will have to know which implementation of GeneMark they are wanting to use for their application, as GeneMark has several variations it can run depending on the desired application. In this work, the GeneMark-ES variant of GeneMark was executed, as it is the self-training *ab initio* GeneMark method for eukaryotic organisms. Options required by GeneMark at runtime are documented in the help message, and simple enough that any user with familiarity of bioinformatics tools should be able to run GeneMark. When running the program, we note that documentation is only provided by the help message, and we were unable to find any additional resources online. In regard to run-time, running the GeneMark-ES pipeline on DC1 with 56 threads finished in 16 minutes. Applying the GeneMark prediction pipeline to Tsth20 took a similar amount of time. GeneMark’s outputs consist of a GTF or GFF file of predicted genes as well as a number of other outputs related to the run.

Braker2 [19] is hosted on GitHub as a repository that receives relatively frequent updates with Braker3 being released before the completion of this work. Braker is maintained by Katharina Hoff from the University of Greifswald and is available under the Open Source Artistic License. Installation of the repository is a straightforward pull from GitHub. As with GeneMark, Braker2 uses a combination of Perl, Python and other executables in its regular use. Downloading the repository itself is not enough for execution, as Braker2 relies on a number of dependencies and bioinformatics tools including Perl and (Perl dependencies), Augustus [48], BamTools [49], BedTools [50], GeneMark [15], StringTie [51], GFFRead [52] and several others. Manual installation of these dependencies would be difficult, time-consuming and in general advised against. In this case, many of Braker2’s requirements are satisfied by modules already included in the environment, making installation relatively simple if one is familiar with the Digital Research Alliance of Canada’s software

stack. This case still required installation of some Perl modules in addition to loading necessary modules. Installation of the Perl modules is straightforward, with documentation available online, and loading modules is also simple, but may be less familiar to users that are unfamiliar with HPC environments. Alternatively, one could use a package manager like Conda to handle installation of Perl modules and other dependencies. This is perfectly reasonable, but also requires knowledge specific to Conda, which can be complicated and frustrating for users with little software management experience. Installation using Conda may also result in duplicate copies of software being installed, wasting disk space and potentially causing conflicts. Once installation is finished, the Braker2 pipeline is relatively straightforward to run as well, with excellent documentation included both online and through the help message. The training, including RNAseq alignments, and gene finding pipeline in this case took 1 hour and 17 minutes using 60 threads. Applying the Braker2 trained gene finding model to DC1 with 60 threads took 21 minutes to complete. Applying the Braker2 pipeline to Tsth20 took a similar amount of time. Run times will of course vary depending on processing power available to the end user, but in the case of *Trichoderma* genomes, users can expect quick results. Once annotation is complete, Braker2 produces a GFF file containing predicted genes along with CDS sequences and amino acid sequences their protein products.

In summation, Braker2 and GeneMark are not direct plug-and-play software packages. Users should expect to encounter issues when getting these programs running in addition to normal downloading and unpacking of software packages, so some expertise is recommended. Neither Braker2 nor GeneMark require users to compile software, however Braker2’s dependencies may require additional compilation and attention. Once installed, both tools are relatively simple to use with documentation available for both on the command-line and excellent documentation available for Braker2 on their GitHub page. In these results, GeneMark runs faster than Braker2, and for Braker2, the training step took roughly three times as long as the gene prediction step. Outputs from both tools are similar although Braker2 has the ability to output coding and amino acid sequences for downstream processing. Both tools run in reasonable amounts of time, where in the case of smaller genomes such as DC1 and Tsth20, users can expect results within a few hours to a day depending on number of computing cycles available to them.

5.4 Initial Gene Finding Results

Prior to discussing gene finding results, we will first define the terms gene and coding sequence in the context of this work. We refer to a gene as a set of start and stop coordinates in a genomic sequence. This definition of a gene simplifies processing, although the definition could be expanded in the future to include introns, exons and potential up and downstream sequences as well. Gene finders may also predict isoforms, or alternative splice variants of a gene based on evidence provided in training, resulting in multiple potential coding sequences for a gene. In this work, we refer to coding sequences as the set of all coding sequences predicted by a gene finder for a given input genome.

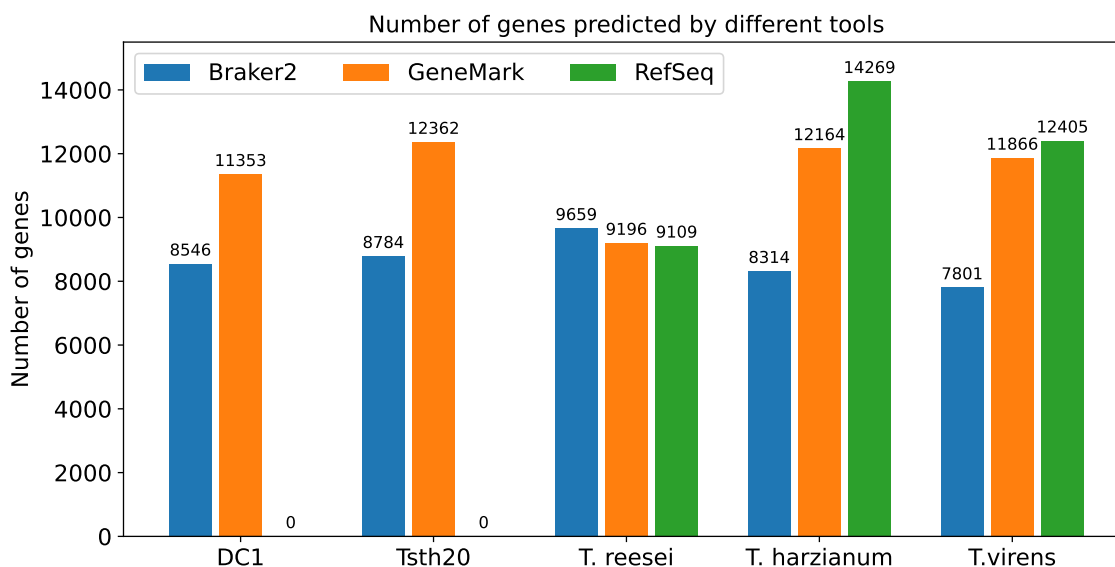


Figure 5.2: Number of genes predicted by each gene finder for each *Trichoderma* genome assembly.

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

Table 5.2: Number of genes and coding sequences (isoforms) predicted by each gene finder for each *Trichoderma* genome.

Counts of genes and coding sequences predicted by Braker2, GeneMark and RefSeq are shown in Figure 5.2, Figure 5.3 and Table 5.2. Immediately we see that Braker2 predicts far fewer genes in all assemblies, except in the case of *Trichoderma reesei*. This is possibly due to the effects of training the Braker2 gene model using data from *Trichoderma reesei*, which has a significantly smaller genome in comparison to other *Trichoderma* assemblies, although genome size is not always indicative of gene content. Regardless, we observe a difference in the number of genes predicted by Braker2 in comparison to GeneMark and RefSeq. The number of genes predicted by GeneMark and RefSeq are similar, except in the case of *T. harzianum*, in which RefSeq predicts roughly 17% more genes than GeneMark. Braker2 consistently predicts more coding sequences than GeneMark and RefSeq. RefSeq also appears to predict multiple coding sequences for each gene but in fewer numbers than Braker2. Coding sequence prediction counts in *T. harzianum* from RefSeq are also interesting, with RefSeq predicting fewer coding sequences than genes. Why this occurs is unknown but may warrant further investigation.

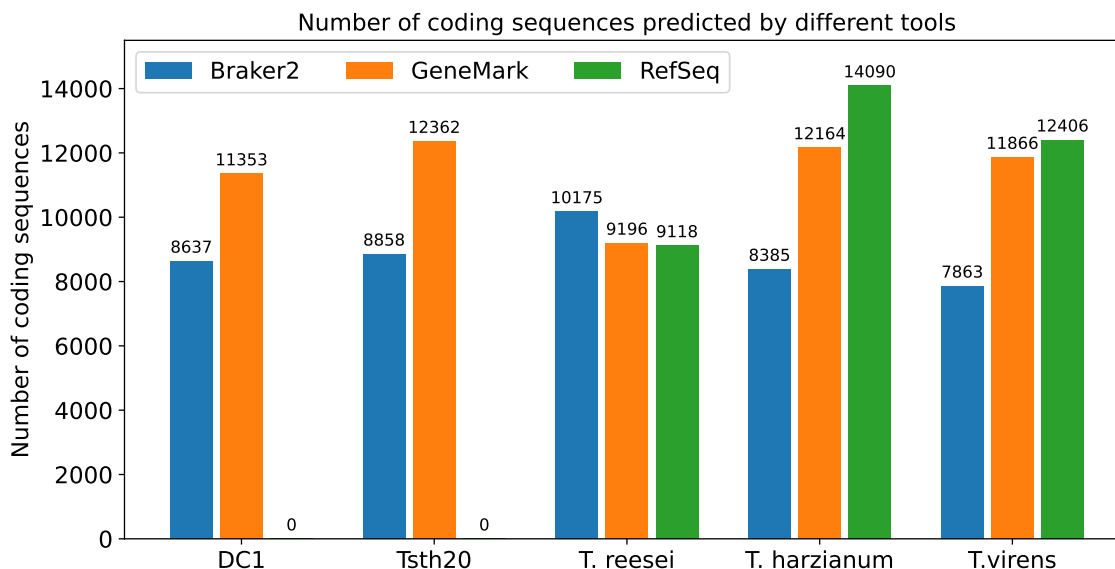
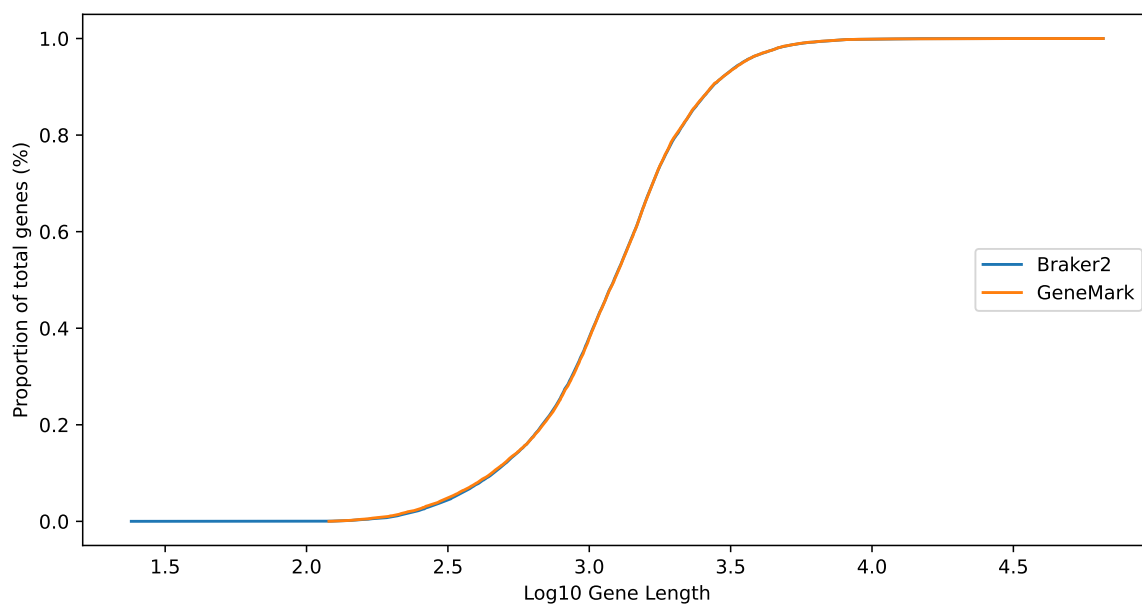


Figure 5.3: Number of coding sequences predicted by each gene finder for each *Trichoderma* genome assembly.

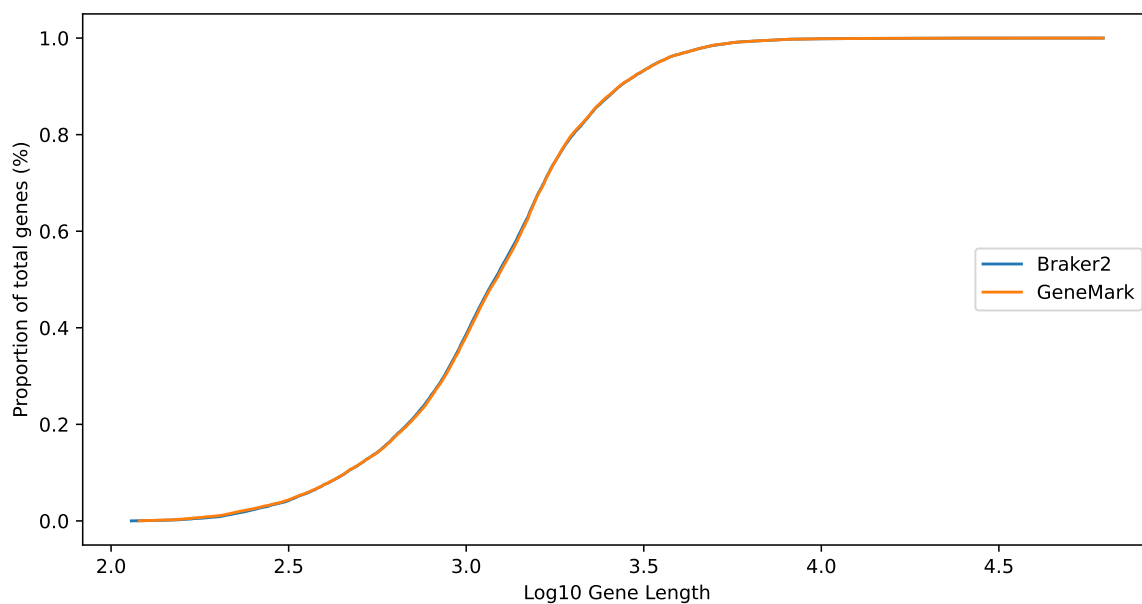
5.5 Distribution of Predicted Coding Sequence Lengths

To better understand and compare the distributions of CDS sequences predicted by Braker2, GeneMark, and RefSeq, the cumulative distribution function for the lengths of CDS sequences for each gene finding tool are shown in Figures 5.4, 5.5 and 5.6. The \log_{10} values of gene lengths were used as the abscissa for a better visualization of the distributions. In DC1, the curves from Braker2 and GeneMark follow each other closely, with the only variation being genes of short length, where Braker2’s curve extends beyond that of GeneMark, indicating that Braker2 predicts the shortest genes in all assemblies except *T. virens*. In the case of Tsth20, the curves are nearly identical. In *T. reesei*, we see disagreement in the curves for shorter genes, with Braker2 appearing to predict a larger fraction of shorter genes than GeneMark and RefSeq. The right sides of the curves trend toward similar predicted gene lengths. In *T. harzianum*, RefSeq deviates from GeneMark and Braker2, predicting more genes of short length, while Braker2 and GeneMark appear to be in near-complete agreement except in the case of very short genes. Finally, in *T. virens*, we see that RefSeq predicts a larger fraction of shorter genes once again, although the deviation is not as drastic as in *T. harzianum*. From these plots, we can say that visually, gene finding tools appear to predict different lengths of genes. We also observe that Braker2 typically predicts the shortest genes of all three gene finding methods.

To confirm that CDF curves differ, two-sided two-sample Kolmogorov-Smirnov[38] tests were performed using the \log_{10} transformed gene lengths, with the null hypothesis being that the difference between the distributions of frequencies of gene lengths predicted by any two tools can be explained by random chance. Results are presented in Table 5.3. In the cases of DC1 and Tsth20, we see that in agreement with Figure 5.4, Braker2 and GeneMark do not produce statistically different lengths of genes, which is interesting considering

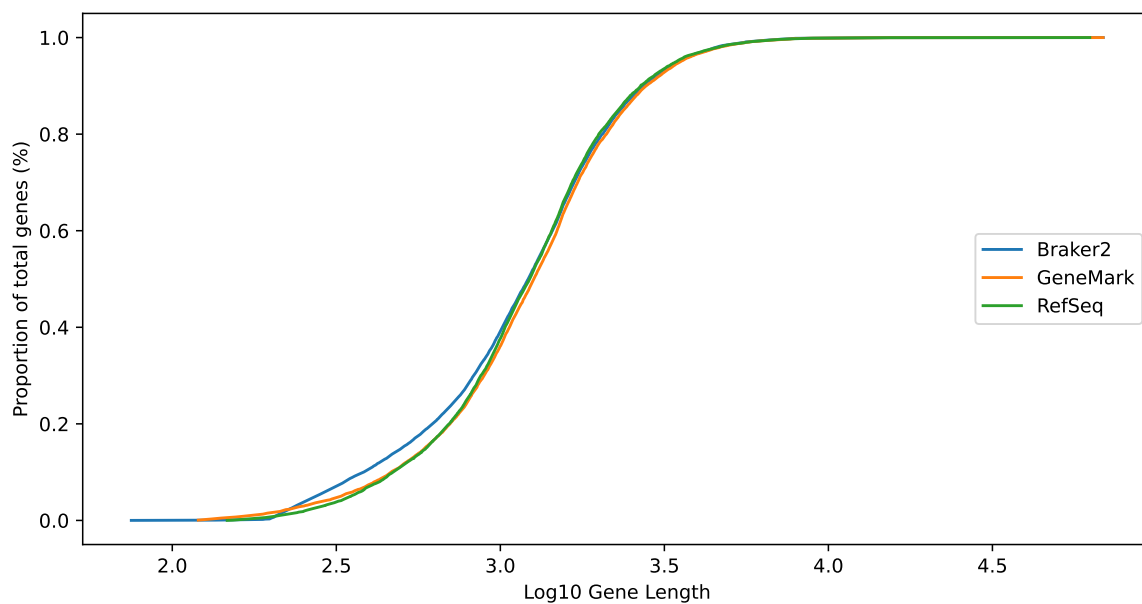


(a) DC1

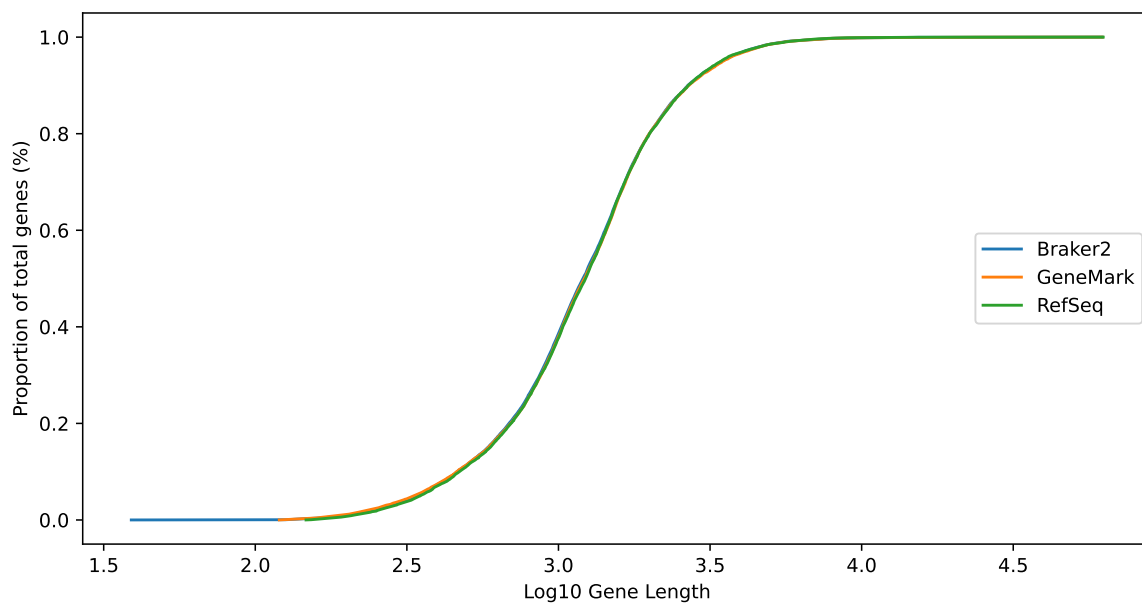


(b) Tsth20

Figure 5.4: Plots of the cumulative distribution function for CDS lengths produced by each gene finding tool applied to DC1 (a) and Tsth20 (b)

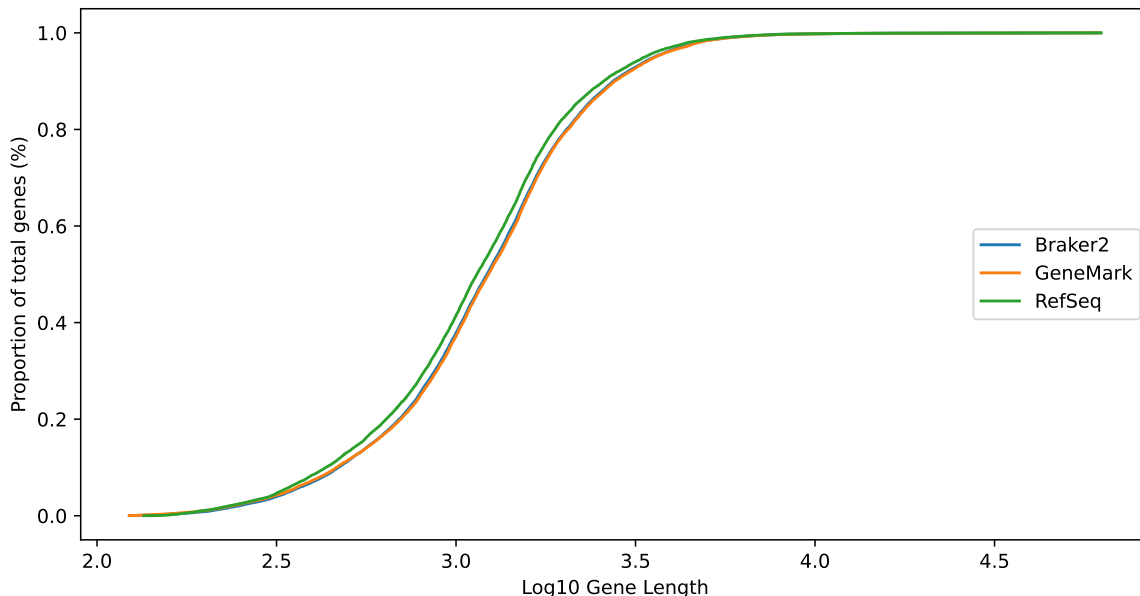


(a) *T. reesei*



(b) *T. harzianum*

Figure 5.5: Plots of the cumulative distribution function for CDS lengths produced by each gene finding tool applied to *T. reesei* (GCF_000167675.1_v2.0)(a) and *T. harzianum*. (GCF_003025095.1_Triha_v1.0)(b)



(a) *T. virens*

Figure 5.6: Plots of the cumulative distribution function for CDS lengths predicted by each gene finding tool when applied to *T. virens* (GCF_000170995.1_TRIVI.v2.0).

that the Braker2 includes experimental evidence for another assembly while GeneMark does not. In *T. reesei*, Braker2’s predicted gene lengths are significantly different from both RefSeq and GeneMark, which is also evident in the CDF plots. The same cannot be said for *T. harzianum* and *T. virens*, where RefSeq is significantly different from both GeneMark and Braker2, which are not significantly different from each other. It is also notable that RefSeq and GeneMark predict similar gene lengths in *T. reesei*, but not in *T. harzianum* and *T. virens*.

It can clearly be stated that these gene finding tools predict different distributions of gene lengths, particularly in *T. reesei*, *T. harzianum* and *T. virens*. Why that may be the case is difficult to answer, but warrants further investigation. There is clearly an underlying difference between RefSeq and the other gene finding tools. Braker2 and GeneMark tend to be in agreement, except in the case of *T. reesei*, for which Braker2 was specifically trained. We observe that when Braker2 is applied to an assembly from which the training data originated, Braker2 predicts a larger fraction of shorter genes. Conversely, we observe that when Braker2 is trained with experimental evidence and applied to a *Trichoderma* assembly from which the training data did not originate, the distributions of predicted genes lengths do not significantly differ from the *ab initio* gene finder GeneMark.

Genome	Tool #1	Tool #2	<i>P</i> -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	$1.794 * 10^{-09}$

Table 5.3: Table of *P*-values from two-sided two-sample Kolmogorov-Smirnov tests between gene finding tools.

5.6 BUSCO Results

Results of BUSCO analysis using the sordariomycetes_Odb12 dataset provided by BUSCO are presented in Figure 5.7, and Table 5.4. The results indicate that all gene sets considered in this analysis have a BUSCO completeness of 94.1% or higher, with a maximum completeness of 98.4% in the case of Braker2 and DC1. In general, Braker2 and RefSeq have the most BUSCO complete sets of gene predictions of the three tools considered. Interestingly, Braker2 produces far more duplicated BUSCO matches than both GeneMark and RefSeq. Examining the BUSCO output logs, this appears to be due to Braker2 predicting more than one coding sequence for some genes predictions, resulting in multiple similar proteins. Interestingly, the coding sequences in the RefSeq annotations seem to miss more genes than the other two gene finders while also having a higher number of fragmented BUSCO genes. This may be due to human curation of the RefSeq datasets, or the Gnomon gene prediction pipeline used by NCBI to produce these annotations. Further investigation is required to determine the exact cause of this discrepancy. Finally, it appears that *T. reesei* tends to have slightly lower BUSCO completeness than the other *Trichoderma* species considered in this analysis, regardless of gene finder used. Why this is the case is unknown, but may be due to the evolutionary distance between *T. reesei*, the Gnomon annotation process, or potentially the fragmented nature of the *T. reesei* assembly used in this analysis. While these results do not capture the entire set of genes possibly present in these *Trichoderma* assemblies, they do confirm that the gene finders are at minimum predicting many evolutionarily conserved fungal genes.

While BUSCO matches are a good metric for general performance of gene finders, it is also important to

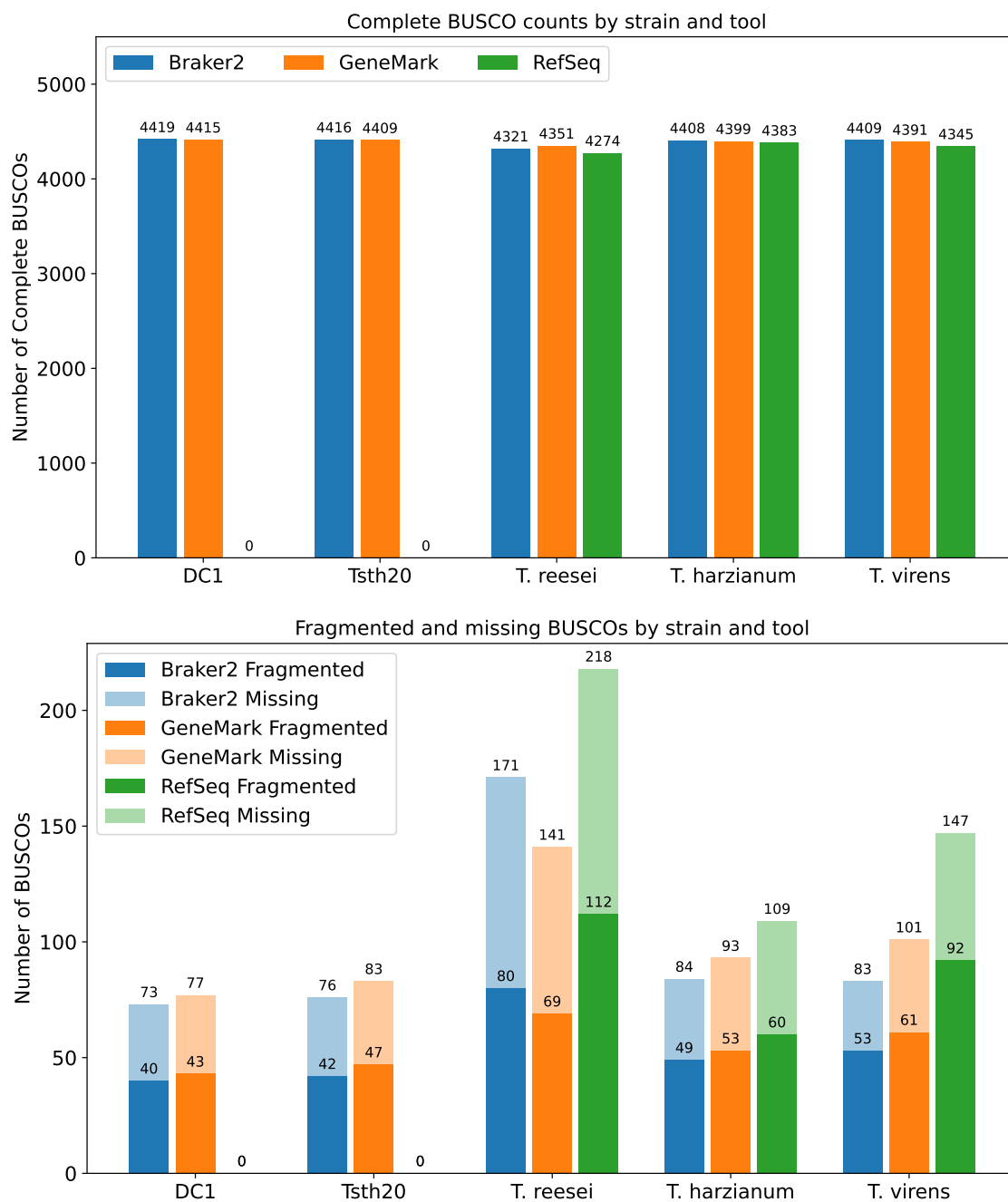


Figure 5.7: BUSCO complete and missing gene counts for each gene finder across all *Trichoderma* genome assemblies. There were a total of 4492 markers in the sordariomycetes_Odb12 BUSCO dataset. For DC1 and Tsth20, RefSeq annotations are not available, so their values are set to 0.

Strain	Complete	Single	Duplicated	Fragmented	Missing
DC1	4419	3547	872	40	33
Tsth20	4416	3585	831	42	34
<i>T. reesei</i>	4321	3587	734	80	91
<i>T. harzianum</i>	4408	3572	836	49	35
<i>T. virens</i>	4409	3530	879	53	30

(a) Braker2

Strain	Complete	Single	Duplicated	Fragmented	Missing
DC1	4415	4401	14	43	34
Tsth20	4409	4392	17	47	36
<i>T. reesei</i>	4351	4345	6	69	72
<i>T. harzianum</i>	4399	4382	17	53	40
<i>T. virens</i>	4391	4369	22	61	40

(b) GeneMark

Strain	Complete	Single	Duplicated	Fragmented	Missing
<i>T. reesei</i>	4274	4267	7	112	106
<i>T. harzianum</i>	4383	4366	17	60	49
<i>T. virens</i>	4345	4321	24	92	55

(c) RefSeq

Table 5.4: Results from BUSCO using the fungal analysis option organized by gene finding tool. The sordariomycetes_Odb12 dataset contains 4492 markers. For more information on the categories assigned by BUSCO, please refer to the BUSCO user guide.

investigate BUSCO proteins that were missing from the gene predictions. Of interest are BUSCO proteins that are systematically missed by a gene finder in multiple or all assemblies, and whether these BUSCO proteins are also missed by other gene finders. Table 5.5 lists 17 BUSCO proteins and their corresponding annotations that were not found in any set of predictions from Braker2, GeneMark or RefSeq. The annotations of these BUSCO proteins do not indicate any obvious reason why they would be systematically missed by all three gene finders, and further investigation is required to determine why these genes are not being predicted. It is possible that these genes are simply not present in these *Trichoderma* species, but this may be unlikely given the evolutionary conservation of BUSCO proteins. Each gene finder also has one or more BUSCO proteins that are missed in all assemblies but not by one or both of the other gene finders. These genes are presented in Table 5.6. Again, the annotations of these BUSCO proteins do not indicate any obvious reason why they would be systematically missed by a particular gene finder, and further investigation is required to determine why these genes are not being predicted. While not presented here, we also identified a number of other genes that are missed by the gene finders which may provide insight into their performance and the evolution of *Trichoderma*. Overall, it appears that there are very few BUSCO proteins that are systematically missed by all gene finders, indicating that the gene finders are generally capable of predicting the majority of conserved fungal genes.

Braker2, GeneMark and RefSeq all demonstrate excellent coverage of the BUSCO fungal protein set, indicating that these gene finders are capable of predicting genes that are expected to be present in these assemblies. From this we can say that the foundations of the underlying gene models used by each gene finder are solid. Braker2 produces more duplicate matches than GeneMark and RefSeq, but this is likely due to multiple isoforms of possible genes being present in the input data. Despite excellent coverage of the BUSCO fungal proteins, all three gene finders miss some BUSCO proteins in their predictions. Finally, we reiterate the possibility that human curation of RefSeq datasets is responsible for these differences, but this requires further investigation.

5.7 Identifying Secondary Metabolite Biosynthetic Gene Clusters

Due to time constraints, only a cursory analysis of secondary metabolite biosynthetic gene clusters was performed on the Tsth20 assembly using the software *antiSMASH*. Tsth20 was selected for this analysis as it is the most contiguous assembly produced in this study, and it is of particular interest due to its observed interactions with plants.

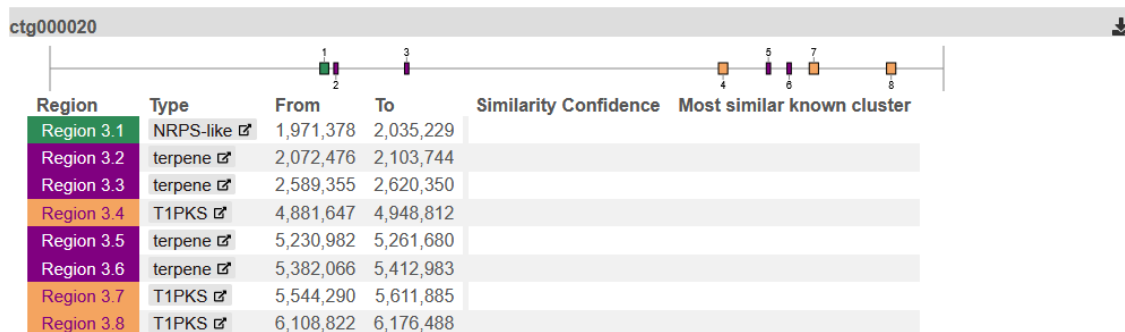
From the Braker2 gene predictions, antiSMASH identified 35 candidate secondary metabolite biosynthetic gene clusters in the Tsth20 assembly, while antiSMASH identified 58 clusters from GeneMark predictions. These clusters varied in length from 30kb to 130kb with no clear separation in size between clusters identified from Braker2 and GeneMark predictions. The RefSeq annotation for Tsth20 was not available for analysis with antiSMASH. The types of secondary metabolite clusters identified included non-ribosomal peptide

BUSCO ID	Annotation
191658at147550	Rossmann-fold NAD (+)-binding protein
250641at147550	Aspartic-type endopeptidase
279527at147550	Phosphatidic acid-preferring phospholipase A1, contains DDHD domain
578862at147550	Transcription factor
579967at147550	Alpha-L-rhamnosidase C
627082at147550	Ferric reductase
628206at147550	Zinc finger domain-containing protein
636196at147550	ATP-dependent RNA helicase
646880at147550	Lipase class 3
647592at147550	Conserved hypothetical protein
652375at147550	Conserved hypothetical protein
657983at147550	Conserved hypothetical protein
658912at147550	Aromatic amino acid aminotransferase
659938at147550	HAUS augmin-like complex subunit 1
672116at147550	Nitrogen regulatory protein AreA, GATA-like domain
677025at147550	Mating-type switching protein Swi10
689133at147550	Myosin heavy chain

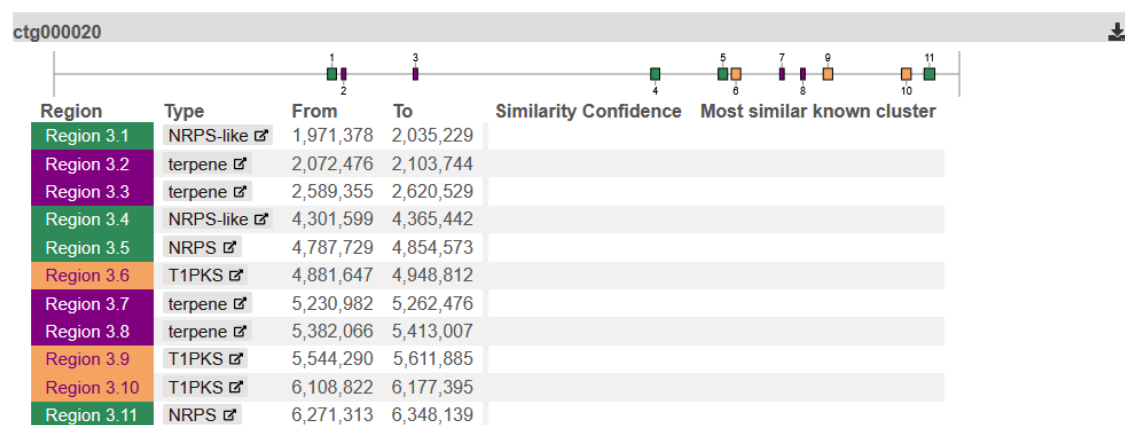
Table 5.5: BUSCO IDs and their corresponding annotations that were not found in any set of predictions from Braker2, GeneMark or RefSeq.

synthetases (NRPS), polyketide synthases (PKS), and terpenes. Example clusters identified by antiSMASH from Braker2 and GeneMark gene predictions from contig ctg000020 are shown in Figure 5.8. We note that the columns addressing similarity to known clusters are left empty, as antiSMASH was run without the option to include this information. From Figure 5.8, we can see that similar clusters were identified by both gene prediction methods, although more clusters were identified from the GeneMark predictions. This may be a result of Braker2 being trained on *T. reesei* gene models, which could lead to a bias in predictions towards genes similar to those in *T. reesei* resulting in missed or incorrectly predicted clusters. The trend of GeneMark gene predictions being more sensitive than Braker2 predictions continues for the other contigs in Tsth20. We also observe that antiSMASH identified no clusters from Braker2 predictions in ctg000000, while 6 clusters were identified from GeneMark predictions in the same contig.

While these results are promising, a more thorough analysis of secondary metabolite gene clusters across all assemblies using multiple gene prediction methods would be necessary to draw any meaningful conclusions about the secondary metabolite potential of these *Trichoderma* strains. Future work could also involve cross-referencing identified clusters with known secondary metabolites and functional annotations to better understand their roles and biosynthetic pathways.



(a) Clusters from Braker2 gene predictions



(b) Clusters from GeneMark gene predictions

Figure 5.8: Example secondary metabolite gene clusters identified by antiSMASH from Braker2 (top) and GeneMark (bottom) gene predictions from contig ctg000020 in the Tsth20 assembly. The columns addressing similarity to known clusters are left empty, as antiSMASH was run without the option to include this information.

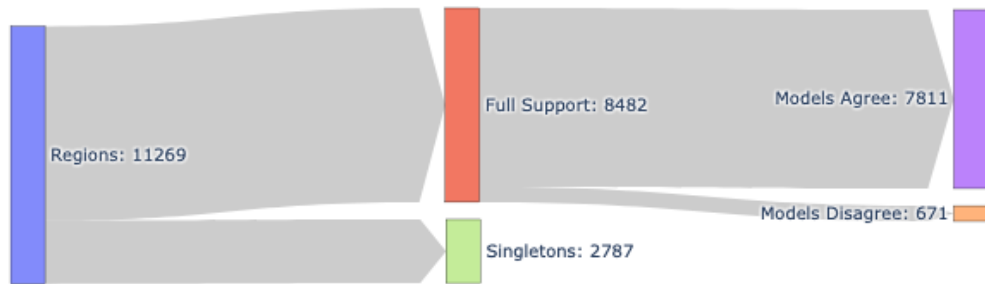
BUSCO ID	Tool	Annotation
282444at147550	RefSeq	Conserved hypothetical protein
291262at147550	RefSeq	HAUS augmin-like complex subunit 6, N-terminal
632369at147550	RefSeq	Peptidyl-prolyl cis-trans isomerase, FKBP-type
632579at147550	RefSeq	Cyanate hydratase
672871at147550	RefSeq	MAU2 chromatid cohesion factor
677279at147550	RefSeq	Pentatricopeptide repeat domain-containing
688724at147550	GeneMark &Braker2	Putative protein of unknown function

Table 5.6: Additional BUSCO IDs along with their corresponding annotations not found in any set of predictions by one or more tools but not all. The tools that missed each BUSCO ID are also listed.

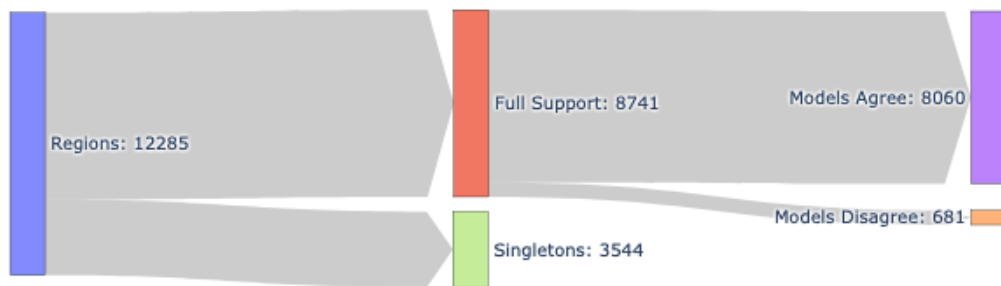
5.8 Region Identification

Visual breakdowns of the types of regions identified and their counts for each *Trichoderma* assembly are shown in Figure 5.9. We will discuss the types of regions, beginning with fully supported regions. These regions, labelled ‘Full Support’, are sections of genomic sequence where all three gene finders agree that a gene is present in some form. These fully supported regions are then broken down into two categories, based on whether the models from each gene finder agree on the start and/or stop positions of the gene model. Regions that have support from more than one gene finder, but not all, are labelled regions with ‘Partial Support’. These regions are also broken down into two subtypes based on whether the gene predictions agree on the start and/or stop positions of the gene. Regions with support from only one gene finder are labelled as singletons.

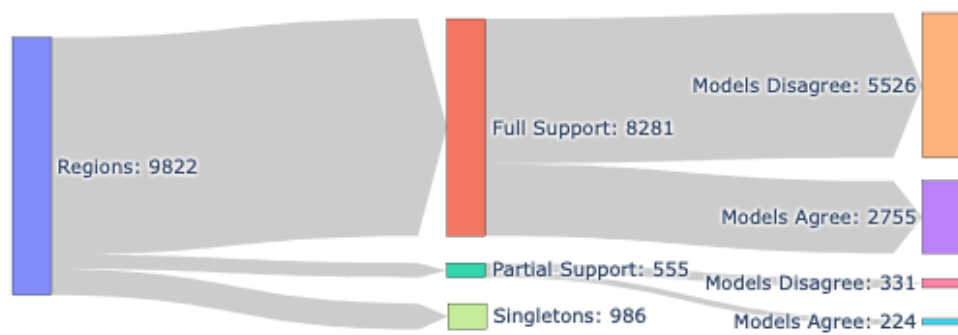
Looking at DC1 and Tsth20 in Figure 5.9, it appears that Braker2 and GeneMark predict genes in the same regions in 69% and 65% of cases, respectively. For regions with full support in DC1 and Tsth20, the models also tend to agree on the start and stop positions of the gene in that region. This is not generally the case as we will see later when RefSeq is also considered. There are no regions with partial support for DC1 and Tsth20, as there are only two gene finders considered. *Trichoderma reesei* contains the fewest number of regions, which seems to scale appropriately with the total number of predicted genes and assembly length. Eighty-four percent of the regions are fully or partially supported by Braker2, GeneMark and RefSeq, with only 10 percent of the regions being single tool predictions. The most interesting observation here is the extent to which fully supported regions disagree on start and stop positions of the gene(s) in each region. There is clearly a difference in the gene models being produced by these gene finders. Also, if there is more disagreement on the start and stop positions of a gene, then there is likely disagreement on the number and location of exons and introns within the gene model as well. In the case of partially supported regions, there is roughly 50% to 60% agreement on start and stop positions, but that may come as less of a surprise as there is already disagreement on presence by definition. Gene finding behavior in regions from *T.harzianum* and



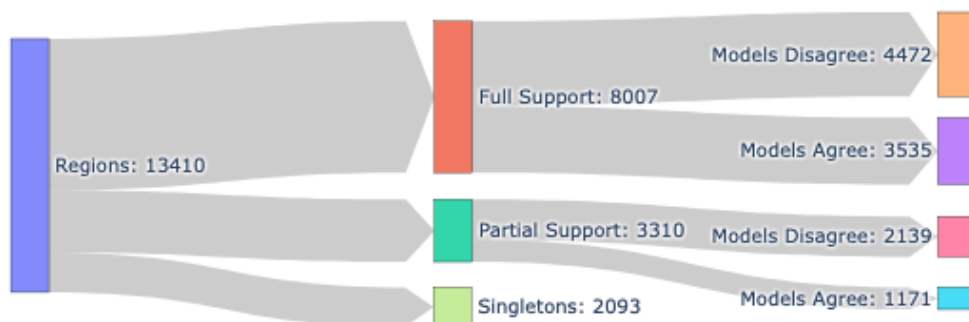
(a) DC1



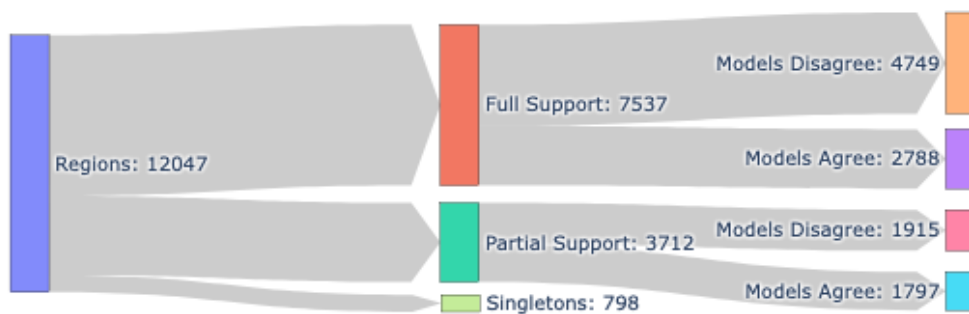
(b) Tsth20



(c) *T. reesei*



(d) *T. harzianum*



(e) *T. virens*

Figure 5.9: Figures showing breakdowns of genomic regions identified by the region finding process. Regions, in blue, are categorized based on support from gene finders. Regions with supporting predictions from all gene finders included in this analysis are labelled with ‘Full Support’ and colored red. Fully supported regions are then broken down into regions where gene models agree on the start and stop positions of the genes (purple), and those that do not (orange). Regions labelled with ‘Partial Support’ (turquoise) are regions with supporting gene predictions from two gene finders. Partially supported regions are also broken down into regions in which gene finders agree on the start and stop positions of the gene (cyan), and those that do not agree (pink). Regions with gene predictions from only one gene finder are labelled as singletons (green).

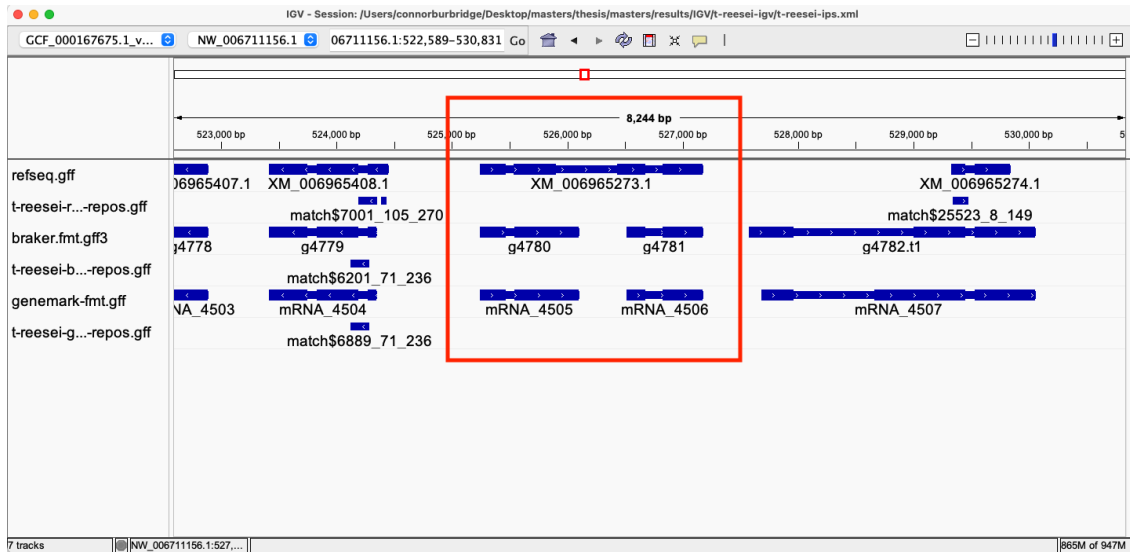


Figure 5.10: An IGV screenshot from *T. reesei* showing a region containing five gene predictions in which one of the gene finders disagrees with the other two, highlighted in the red box.

T. virens differ from the other assemblies in the split between fully and partially supported regions. There seems to be fewer regions with full support from all gene finders and the reason why is unclear. The gene models in each region, whether fully or partially supported, still tend to disagree on start and stop positions of the genes more often than agree.

It is also worthwhile investigating some potentially interesting cases of unusual gene calls in regions. One such case is when a region contains more than three individual predicted genes. The null hypothesis, so to speak, is that if all gene finders agree exactly on the presence and position of a gene in a region, then one should expect exactly three predictions in that region. We observe that this is not always true. Cases of more than three gene predictions in a region were present in regions from all processed assemblies. Figure 5.10 shows an example of one of these unusual regions, in which the single gene predicted by RefSeq spans multiple gene predictions from both Braker2 and GeneMark. In the case of DC1 and Tsth20, there were very few regions with more than three gene predictions, with DC1 having 19 and Tsth20 having 6. Conversely, *T. reesei*, *T. harzianum* and *T. virens* contain many such regions, with assemblies reporting 546, 899 and 521 regions with greater than three gene predictions, respectively. While having predictions from each gene finding tool present in a region is a strong indicator for the presence of a gene, having more than one gene prediction per gene finder raises questions about which model (or models) is correct and why disagreement exists.

Another interesting set of criteria to investigate is whether gene finders always predict genes on the same strand within a region. We observe that regions containing predictions on different strands do exist, and are observed in predictions from all assemblies included in this analysis. Figure 5.11 shows an example of a region containing gene predictions on opposing strands. As in the case of regions with more than three gene predictions, DC1 and Tsth20 have fewer mixed strand regions with 29 regions identified in DC1 and 46

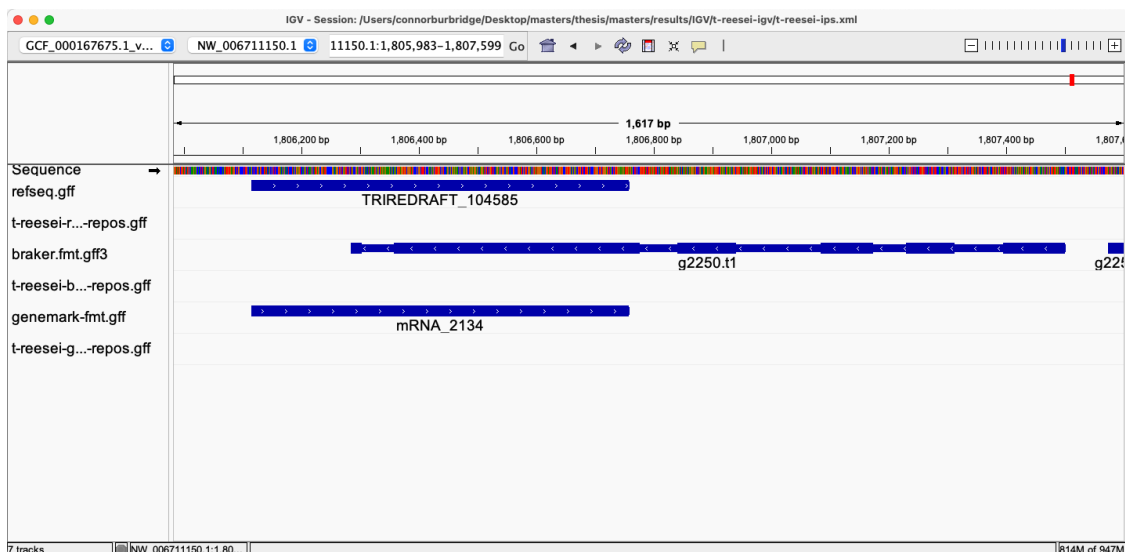


Figure 5.11: An IGV screenshot of *T. reesei* showing a region which contains gene predictions on opposite strands.

regions identified in Tsth20. Results from *T. reesei*, *T. harzianum* and *T. virens* show more regions in which this property is true, with region counts of 203, 533 and 293 respectively. Under the assumption that gene finders should predict the same genes on the same strands, it is unexpected to find so many of these cases, although it is possible that these overlapping sense and anti-sense genes exist naturally [53].

In summary, gene finders agree partially or completely on the presence of a gene in the vast majority of cases. While gene finders generally agree on the presence of a gene, they tend to disagree on the underlying gene model more often than they agree, except in the case of DC1 and Tsth20. This is likely due to only including two gene finders in their analysis rather than three, resulting in fewer opportunities for disagreement. This observation is true when applied to the start and stop positions of the gene, but further investigation and comparison of intronic and exonic sequences between genes may provide more insight. It is also possible that genome quality may affect agreement between gene finders as the underlying model may be trained on contigs that are larger or of different structure than the input genome. This is particularly true in the cases of DC1 and Tsth20, which are composed of larger contigs in comparison to the assemblies from NCBI. Finally, regions identified in this analysis do not always fit the ideal scenario of one gene prediction from each tool per region. Regions in which there are more than three gene predictions were observed in all assemblies included in this work. In addition, we observe cases where gene finders predict genes on different strands.

5.9 InterProScan as Supporting Evidence for Predicted Genes

Pfam hits from InterProScan analysis are presented in Figure 5.12 and Table 5.7, from which we can identify one major trend. In general, roughly 7376% of proteins predicted by Braker2, GeneMark and RefSeq contain a match to a Pfam entry, except in the case of *T. harzianum*, which reports a considerably lower

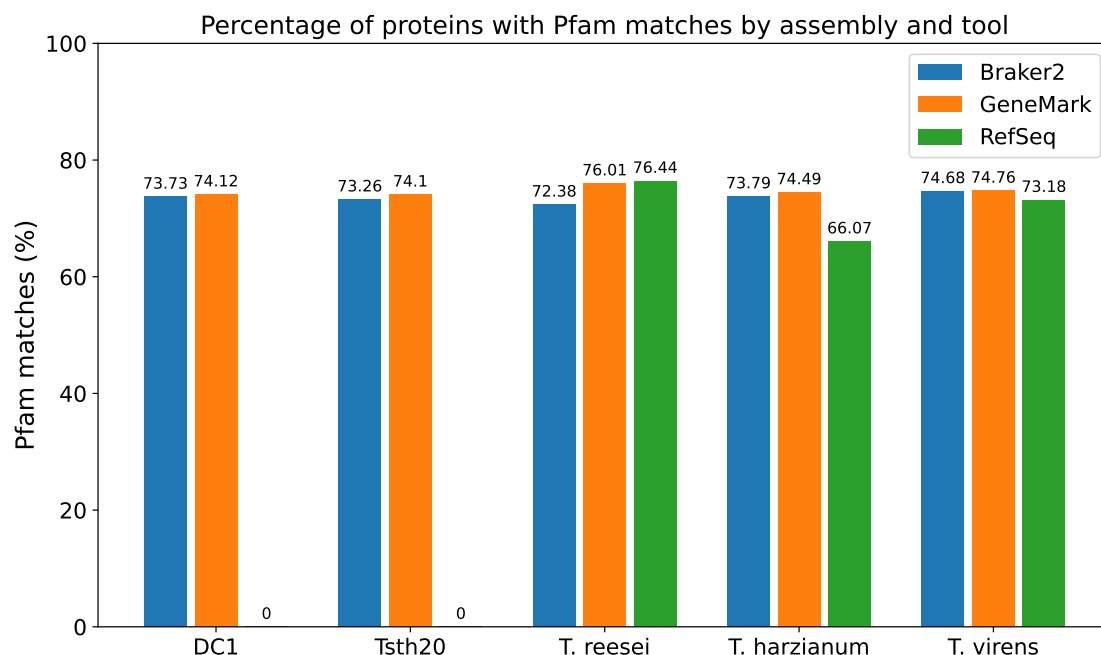


Figure 5.12: Bar plot showing the percentage of predicted proteins with Pfam matches from InterProScan for each gene finder across all *Trichoderma* genome assemblies. For DC1 and Tsth20, RefSeq annotations are not available, so their values are set to 0.

proportion of genes with Pfam matches in the RefSeq dataset. Why the proportion of RefSeq proteins with Pfam matches in *T. harzianum* is so low is unknown. This is promising performance for the gene finders as the RefSeq annotations demonstrate similar proportions of Pfam hits to predicted proteins. The total counts may be deceiving however, as predictions from Braker2 may result in more than one protein product per gene, whereas in the case of GeneMark, only one protein is produced per gene model. From visual inspection of Pfam hits mapped back to the references it appears that, in general, when gene finders agree that a gene is present, InterProScan reports the same Pfam match in all three predictions. An example of agreement between Braker2, GeneMark, RefSeq and InterProScan is shown in Figure 5.13. It is important to note that the position of the Pfam match in IGV does not indicate the true position of the Pfam match in the gene, but provides an indication of the presence of Pfam matches. Pfam matches are offset from the start of the gene based on the start and end position of the Pfam match in the protein sequence. For example, if a Pfam match has a start position 10 amino acids into the protein sequence, the corresponding start position in the resulting GFF is 10bp downstream from the start of the gene.

While cases of complete agreement are abundant, cases of disagreement also exist and in strange forms. In many cases, while the gene finders agree on the presence of a gene, only the RefSeq protein product contains a match to the Pfam database. Why this may be the case is unclear, and may warrant further investigation. There are also many cases in which InterProScan reports the same Pfam hits for individual proteins, but the gene models in the region do not agree. There are even cases such as the region shown in Figure 5.14, where

Assembly	Braker2 (%)	GeneMark (%)	RefSeq (%)
DC1	73.73	74.12	N/A
Tsth20	73.26	74.10	N/A
<i>T. reesei</i>	72.38	76.01	76.44
<i>T. harzianum</i>	73.79	74.49	66.07
<i>T. virens</i>	74.68	74.76	73.18

Table 5.7: Table showing the fractions of predicted genes with Pfam annotations from InterProScan as a percent

Braker2 and GeneMark agree that two genes and their associated proteins and Pfam matches are separate, but RefSeq only reports one gene with multiple Pfam hits. There are also several cases where two tools are in agreement with proteins containing Pfam hits while another is not. Even more interesting are cases such as the one shown in Figure 5.15, in which Braker2 and GeneMark predictions contain Pfam matches while RefSeq does not report a gene at all. This may be due to experimental data used in the RefSeq training process or a result of curation. Regardless of the tool, these cases demonstrate well that gene finders are not always in agreement even on well known proteins, to the point that predictions are not present even though protein products from other gene finders contain known Pfam matches.

In summary, the protein products from Braker2 and GeneMark predictions do contain matches to the Pfam database. The proportions of matches to total proteins are similar to that of the RefSeq annotation, sitting between 65 and 75 percent. While Pfam hits to Braker2, GeneMark and RefSeq proteins generally agree, we observe regions in which there is disagreement in several forms.

5.10 BLAST Results

To determine if gene finders predict genes that are present in other similar species, proteins from three related organisms were used as queries to `tblastn`[54] on each *Trichoderma* assembly. The organisms and their associated NCBI accession IDs are as follows: *T. atroviride* GCF_000171015.1, *F. graminearum* - GCF_000240135.3, *S. cerevisiae* GCF_000146045.2. Results from the `tblastn` runs are presented in Figure 5.16 and Table 5.8. Initial `tblastn` results appear promising for both the *T. atroviride* and *Fusarium* datasets. All assemblies considered contain at minimum 89% of the reference protein sequences in the case of *T. atroviride* and a minimum of 75% in the case of *Fusarium*. In the case of *S. cerevisiae*, a minimum of 57% of reference proteins matched. It appears that the percentage of hits decreases as evolutionary distance increases. These results provide rough validation that the assemblies contain potential for protein coding sequences.

While successful `tblastn` hits from input proteins are useful, it is worthwhile exploring those hits in the context of regions and gene predictions within them. Counts of genes from regions with `tblastn` hits are shown in Figures 5.17a, 5.17b, 5.17c, and Table 5.9. In a similar trend to Table 5.8, genes in regions with `tblastn`

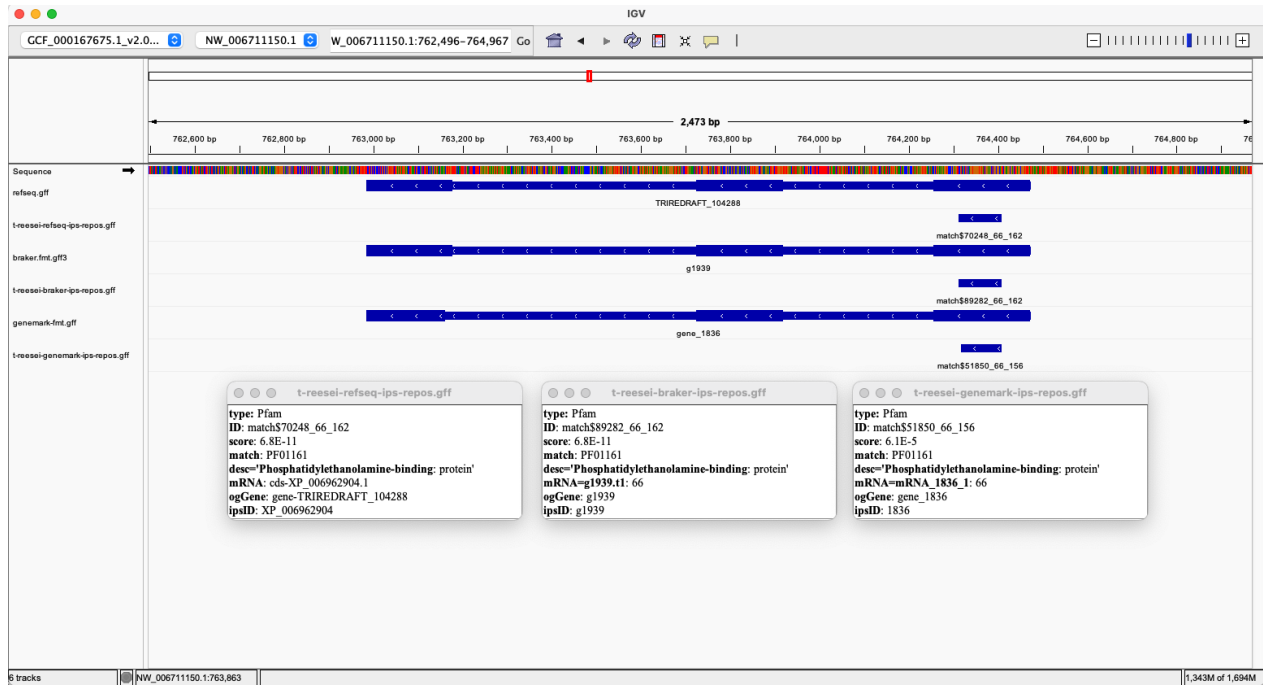


Figure 5.13: An IGV capture showing complete agreement between gene finders for both gene model and protein Pfam hits. The longer segments are predicted genes and the smaller segments are annotated Pfam matches, all with agreeing start and stop positions. The description in each sub-window indicates the agreeing Pfam annotations.

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

Table 5.8: tblastn hits from reference protein sequences to selected *Trichoderma* assemblies. Hits are reported if the alignment length is greater than 30% of the reference protein length and if 30% of the alignment length has identical matches.

hits decrease as evolutionary distance increases. In DC1 and Tsth20, we observe that Braker2 predicts more genes in regions with tblastn hits than GeneMark. RefSeq is not included for DC1 and Tsth20, as there is no RefSeq annotation available. In *T. harzianum* and *T. virens*, Braker2 and RefSeq appear to have similar counts of genes in regions with tblastn hits. Interestingly, in the case of *T. reesei*, there are more genes from GeneMark in regions with tblastn hits, which is not the case in *T. harzianum* and *T. virens*.

Another interesting observation is comparing the coverage, or completeness, of the input and query sequence. As mentioned, it appears that high proportions of query sequences are being mapped to the subject, but the gene predictions are not as well represented in regions where those tblastn hits occur. A likely reason for this is inappropriate alignment parameters. More flexible gap parameters and the use of a different scoring matrix would likely produce better results, but optimization of parameters was not performed in this

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
Tsth20	<i>T. atroviride</i>	6065	4626	N/A
Tsth20	<i>F. graminearum</i>	5365	4191	N/A
Tsth20	<i>S. cerevisiae</i>	2211	1869	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
<i>T. virens</i>	<i>T. atroviride</i>	6256	4437	6415
<i>T. virens</i>	<i>F. graminearum</i>	5568	4075	5664
<i>T. virens</i>	<i>S. cerevisiae</i>	2318	1861	2352

Table 5.9: Counts of transcripts predicted by Braker2, GeneMark and RefSeq in regions containing tblastn hits from *Trichoderma atroviride*, *Fusarium graminearum*, and *Saccharomyces cerevisiae*.

analysis. It is also possible that the post-alignment filtering criteria were inappropriate, resulting in short alignments that are true in nature, but may not be indicative of the presence of a truly similar protein.

Overall, it appears that using tblastn with proteins from related organisms may be a viable option as a control when no reference or model organism is available. At the very minimum, tblastn hits can be used as another piece of supporting evidence for the existence of a gene, with one caveat being imperfect coverage of predictions. The issue of imperfect coverage may be mitigated by parameter optimization. It is also possible that one could employ a similarity-based cutoff on the proportion of predicted genes in regions with tblastn hits as an indicator of gene finding performance; however, this would require further analysis, preferably with more closely related organisms.

5.11 Genes in Regions of Anomalous GC Content

Figure 5.18 and Table 5.10 shows the results of applying the same region finding process from earlier to segments of the genome identified as AT-rich. AT-rich is defined as a region of genomic sequence with percent GC composition less than 28% as determined in Section 5.2. In comparison to results from Section 5.8, we see that overall there are far fewer regions with gene predictions in AT-rich genomic segments. In DC1, Tsth20, there are very few regions with full support from both Braker2 and GeneMark, but more singletons.

Assembly	Full Support	Partial Support	Singletons	No. Genes
DC1	11	N/A	20	42
Tsth20	2	N/A	9	13
<i>T. reesei</i>	25	18	54	194
<i>T. harzianum</i>	26	43	68	265
<i>T. virens</i>	8	11	0	49

Table 5.10: Regions of full and partial agreement as well as singleton regions in AT-rich genomic sequence. It is important to note that in the cases of DC1 and Tsth20, full support indicates supporting gene predictions from both GeneMark and Braker2.

Again, as in Section 5.8, there are no regions with partial support as only two gene finding tools were applied to those assemblies. *T. reesei* and *T. harzianum* report more regions in AT-rich genomic sequence than the other assemblies, but with the majority of regions belonging to the singleton category. *T. virens* is an interesting case, reporting a similar number of regions and genes in AT-rich genomic sequence as DC1 and Tsth20. *T. virens* is also the only assembly to report zero singleton gene predictions. Why *T. virens* differs from the other RefSeq assemblies is unclear. In general, there are few regions with gene predictions in AT-rich genomic segments, and within these regions, the majority of cases are isolated singleton gene predictions. These observations differ greatly from those made in nucleotide composition agnostic approach in Section 5.8.

In addition to a breakdown of regions in AT-rich genomic sequence, understanding which gene finders predict more or fewer genes in these regions may be of interest. It is possible that an HMM, trained on genomic sequence with varying nucleotide content, may predict genes differently to an HMM trained only on sequences with uniformly distributed nucleotide composition as the variation in nucleotide composition may affect the various states and relationships between them in an HMM. Figure 5.19 and Table 5.11 shows the number of genes predicted by each gene finding tool in regions of AT-rich genomic sequence. GeneMark appears to predict the fewest genes in AT-rich regions, while RefSeq appears to predict the most. Braker2 lies somewhere in the middle. Again, *T. virens* appears as an odd case, with very few predictions from all gene finders. Why *T. virens* differs from the other RefSeq assemblies is unclear.

Finally, to test the probability of any given gene prediction falling in an AT-rich genomic sequence, a two-sided binomial test was performed to determine if the number of genes predicted in AT-rich sequences is proportional to fraction of genomic sequence they comprise. The null hypothesis in this case is that the gene finding tools predict the same proportion of genes in AT-rich genomic sequences as they do in typical genomic sequence. The results of the test are shown in Table 5.12. In all cases, it appears that the null hypothesis is rejected.

In summary, very few regions with gene predictions are present in AT-rich genomic sequence. Additionally, genes predicted in these regions tend to be isolated and not supported by other gene finders, although some

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

Table 5.11: Number of genes predicted by Braker2, GeneMark and RefSeq in AT-rich genomic sequence from each assembly.

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	$1.29 * 10^{-49}$	$2.44 * 10^{-205}$	$7.40 * 10^{-33}$

Table 5.12: p -values produced from a two-sided binomial test for each combination of tool and assembly.

agreement is observed. In terms of number of genes predicted, RefSeq tends to predict the most genes in these AT-rich regions while GeneMark predicts the fewest. Lastly, the selected gene finding tools do not predict genes in AT-rich sequences in proportion to the fraction of genomic sequence they comprise.

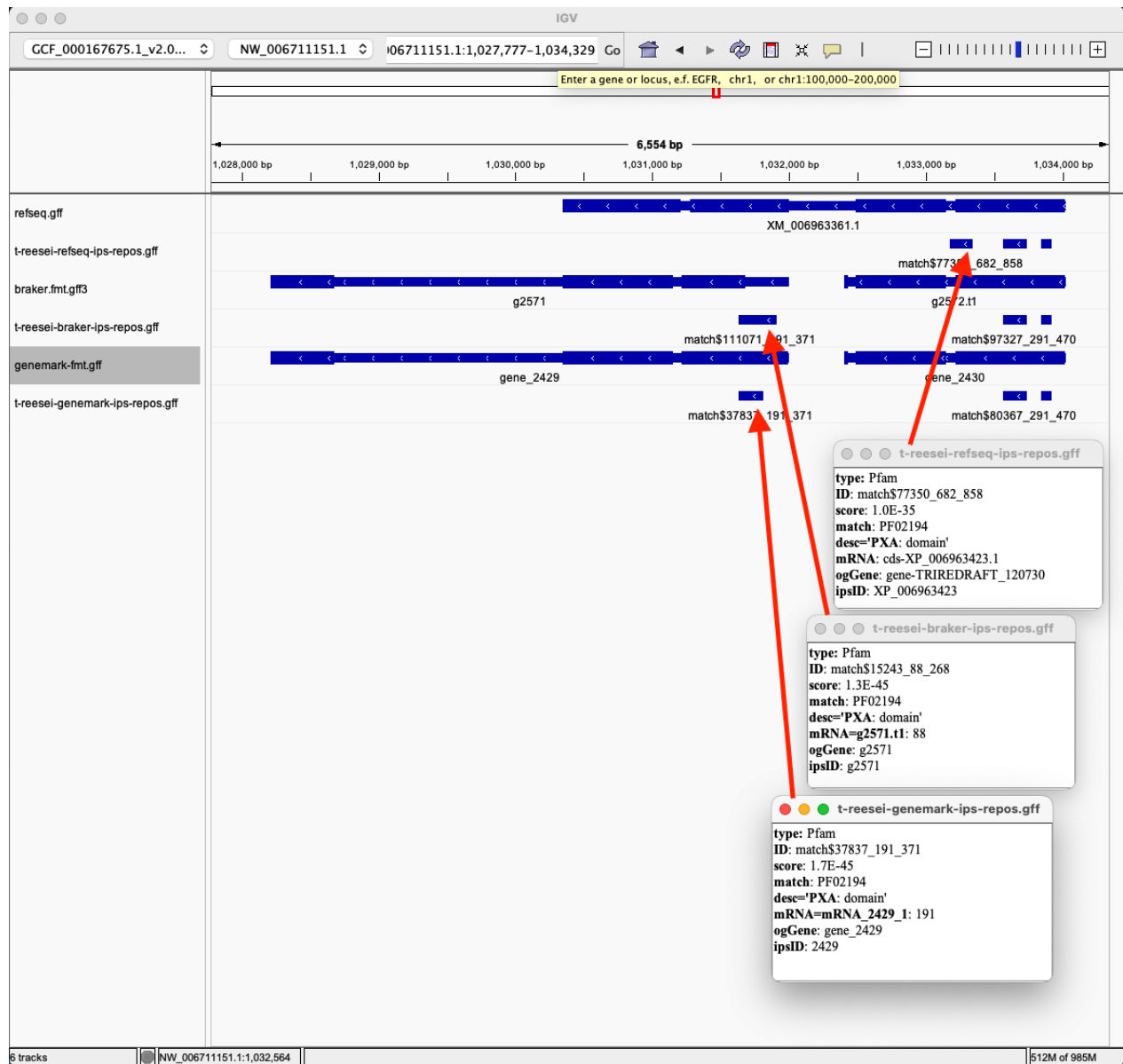


Figure 5.14: An IGV capture showing Braker2 and GeneMark reporting two genes and their resulting proteins and Pfam hits as separate, while RefSeq reports one gene, one protein and three Pfam matches.

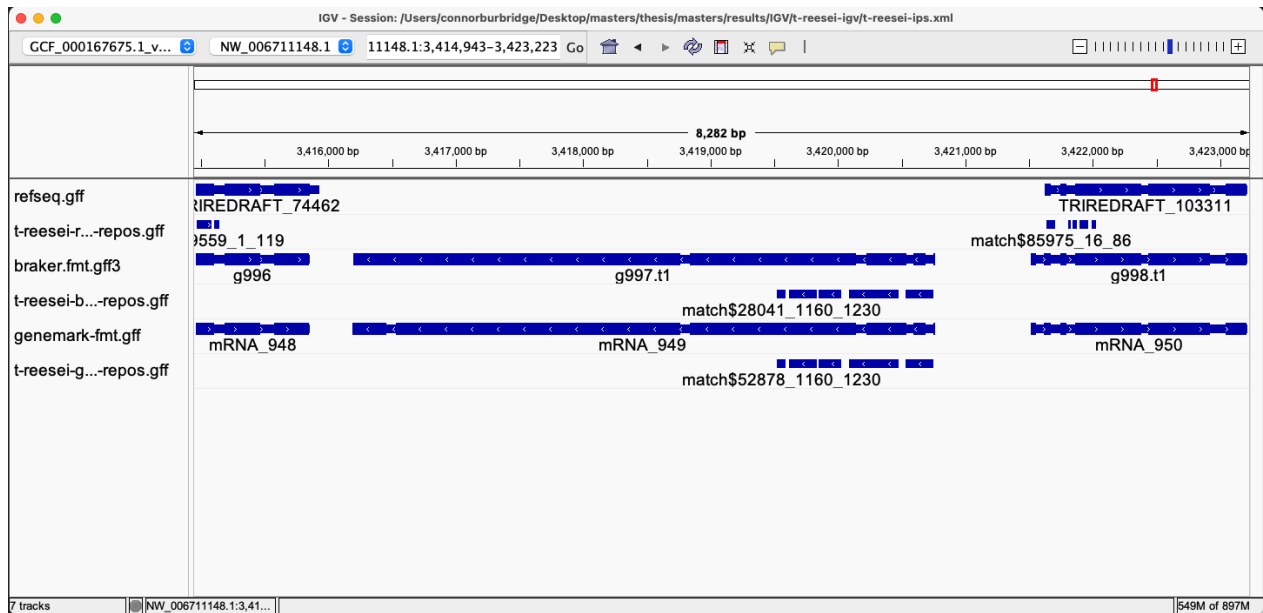


Figure 5.15: An IGV capture showing a scenario where GeneMark and Braker2 agree on a gene model with supporting Pfam evidence and RefSeq does not report any gene.

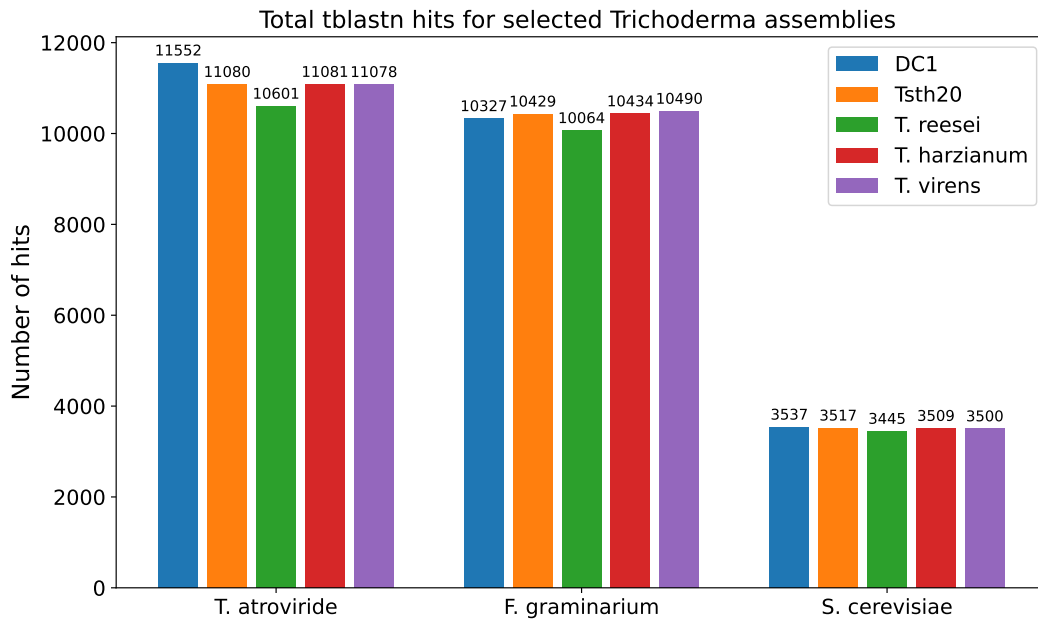
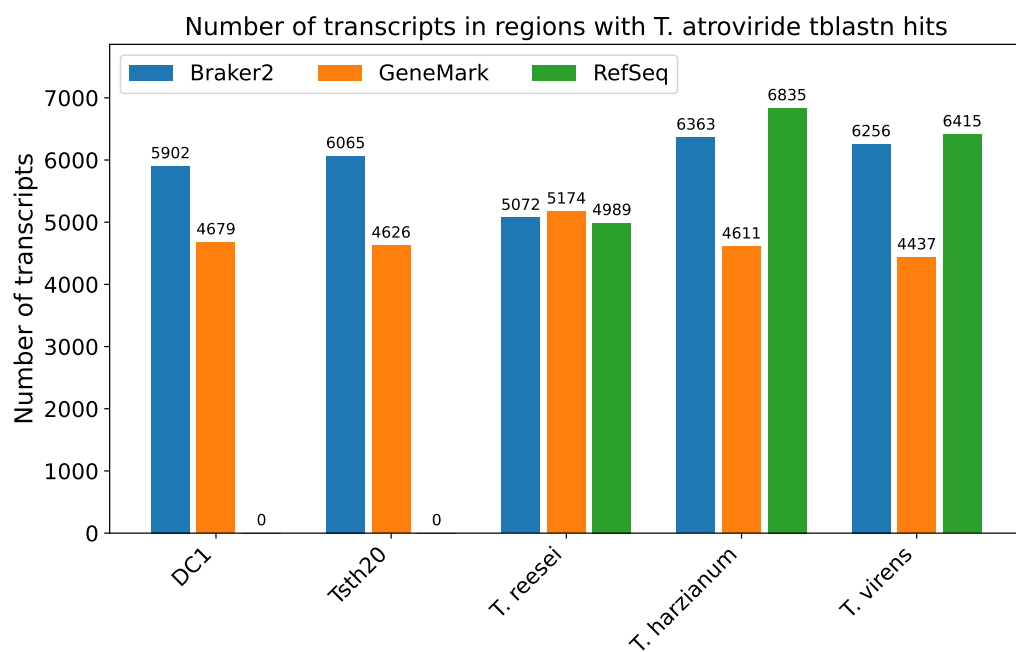
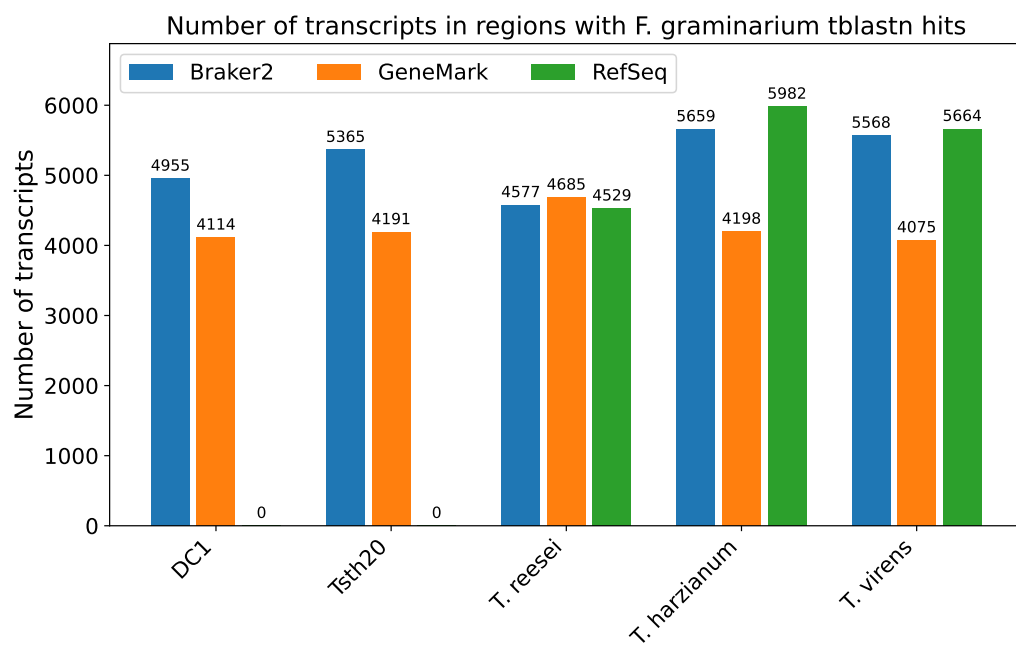


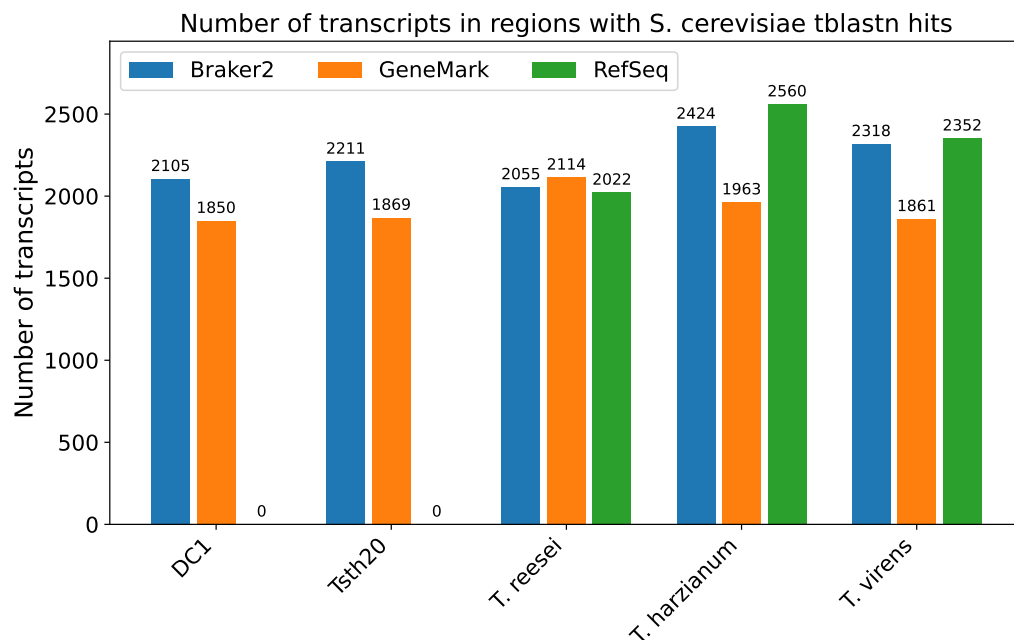
Figure 5.16: Total tblastn hits for reference proteins from *T. atroviride*, *F. graminearum*, and *S. cerevisiae* against *Trichoderma* assemblies.



(a) *T. atroviride* tblastn hits



(b) *F. graminearum* tblastn hits



(c) *S. cerevisiae* tblastn hits

Figure 5.17: tblastn hits for reference proteins from *T. atroviride*, *F. graminearum*, and *S. cerevisiae* against *Trichoderma* assemblies. For DC1 and Tsth20, RefSeq annotations are not available, so their values are set to 0.

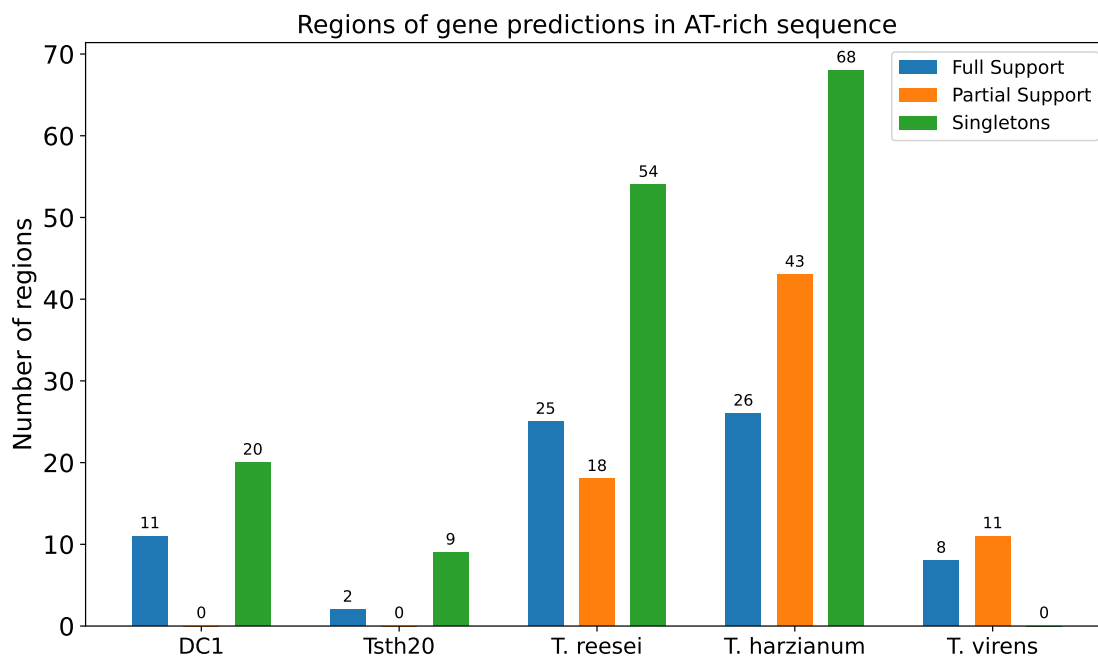


Figure 5.18: Regions of full, partial, and no agreement in AT-rich genomic sequence. It is important to note that in the cases of DC1 and Tsth20, full support indicates supporting gene predictions from both GeneMark and Braker2 and as such, there are no regions with partial support.

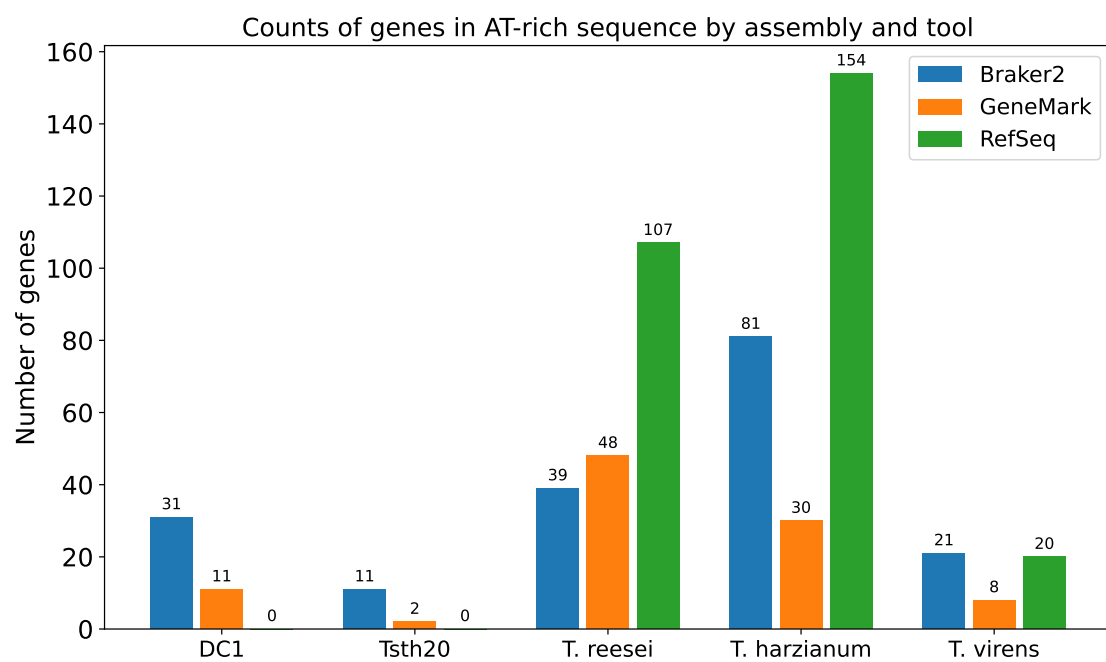


Figure 5.19: Number of genes predicted by each gene finder in AT-rich genomic sequence.

6 Conclusions

6.1 Selection of Gene Finding Tool

With all of these results, it makes sense to explore the question of which gene finding tool one should choose for optimal gene prediction performance. Comparisons drawn in this section are made in the context of *Trichoderma* assemblies and may not extend to other datasets. Observations from the results portion of this work have been converted to scores for each gene finder, relative to the other gene finders. Scores range from zero to three, with values being: 0 failing performance, 1 passing performance, 2 good performance, and 3 excellent performance. Results from this work are summarized in Table 6.1. Ignoring availability and use of the gene finding tools, it would appear that RefSeq performs the best in the remaining categories, earning top marks in every category except in its ability to predict very short genes. Braker2 earns second place; however, this does not capture Braker2’s failure in predicting accurate numbers of genes in DC1, Tsth20, *T. harzianum* and *T. virens*. GeneMark comes in last, excelling only in number of genes predicted and Pfam support for the genes that it predicts. We do note that these criteria are somewhat arbitrary and subjective in some cases, however they do provide a useful framework for comparing the gene finders, and are based on objective results presented in this work.

Relating these observations to use-case scenarios, in the case a researcher’s organism of interest has a RefSeq annotation associated with it, the RefSeq gene prediction process appears to produce the best set of predictions. If users also have experimental evidence, such as RNAseq data under experimental conditions, it may be worthwhile training a Braker2 model and predicting genes with Braker2 to supplement the already well performing RefSeq gene predictions. We note that users should take the quality of the RefSeq assembly into account and investigate the input data and assembly protocols to ensure high-quality predictions. In a similar case, if the organism is of RefSeq status, but the assembly in question is novel to the work being performed, training and prediction using a Braker2 model with experimental evidence either from the novel genome or from the RefSeq accession in addition to use of the RefSeq annotation would likely be the best option. If the organism of interest is not a RefSeq individual but training data is available for that organism, Braker2 is the next best option, although it is important to note that the application of a trained Braker2 prediction model to an organism from which the training data did not originate is not advised based on the results presented in this work (see 5.4). While it is true that the *T. reesei* genome differs from other *Trichoderma* genomes, it’s status as a representative RefSeq organism makes it somewhat of a gold standard. In this case, applying a gene model trained using evidence from the gold standard produces biased numbers

of genes predicted in other *Trichoderma* genomes. While Braker2 technically scores the second highest, users must be very careful when selecting training data, and ensure that the training data either comes from the organism of interest, or comes from a very closely related organism with a highly similar genome. In the case that no appropriate training data is available, GeneMark is still an option, and users can be confident that the tool predicts a reasonably accurate number of genes with supporting Pfam matches. It is also important to note that GeneMark does not perform as well in AT-rich regions as Braker2 and RefSeq, does not predict isoforms, and systematically fails to predict some BUSCO orthologs. Identification of secondary metabolite biosynthetic gene clusters may also be impacted by the choice of gene finder, as GeneMark predicted more genes in the Tsth20 assembly, leading to the identification of more candidate clusters by antiSMASH, however further analysis is required to determine the accuracy of these predictions.

When availability of a gene finding tool becomes a concern and RefSeq is not considered, the scores drop significantly for Braker2 and GeneMark as seen in the final row of Table 6.1. In this situation, Braker2 still outperforms GeneMark. If the organism of interest is not considered a representative RefSeq individual, but supporting evidence specific to that organism or a very closely related organism is available, a trained Braker2 prediction model will perform well. Again, in the case that the organism is not a RefSeq individual and no appropriate training data is available, GeneMark is still a reasonable option even with the previously identified caveats.

In summary, these results indicate that if your organism is a RefSeq organism, use the RefSeq annotation. If no RefSeq predictions are available, but appropriate training data is, one should use Braker2. If the training data is of questionable similarity or not available at all, users can fall back on *ab initio* gene finders such as GeneMark, which while not ideal, still predict genes with supporting evidence.

6.2 Future Work

While this work is extensive, there are still many areas that could be expanded on. The sheer number of gene finding tools available means that many more could be compared to the ones presented here. In addition, supplying more training data to Braker2 may improve gene finding performance, and training on RNAseq from different organisms could be used to improve the performance of Braker2. With more data in mind, different types of training data such as RNAseq, expressed sequence tags (ESTs), and protein sequences could also be used to train models that support external evidence. Another limitation of this work is that the gene finders were run with default parameters, and it is possible that tuning the parameters of the gene finders could improve performance.

Following the gene prediction process, there are a number of different analyses that could be performed on the predicted genes. For example, functional annotation of the predicted genes could be performed using tools such as BLAST2GO to assign Gene Ontology (GO) terms and KEGG pathways to the predicted genes. This could provide further insight into the functional roles of the predicted genes and their potential applications.

Category	Braker2	GeneMark	RefSeq
Availability	3	3	0
Ease of install	1	2	0
Ease of use	3	3	0

# 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

Table 6.1: Table with scores attributed to performance of each gene finder in several categories. The score definitions for performance are as follows: 0 fail, 1 pass, 2 good, 3 excellent. Since RefSeq is not publicly available, it is marked as N/A in the publicly available final scores. The dashed line separates categories associated with operation and use from categories describing gene prediction performance.

More importantly in the context of this work, further analysis of predicted genes with a tool like antiSMASH could be used to identify biosynthetic gene clusters (BGCs) in the predicted genes.

Expanding the number of genomes used in this work would also be beneficial, as the results presented here are based on a limited number of genomes. Including more genomes from different species and strains of *Trichoderma* could provide a more comprehensive understanding of gene finding performance across the genus. Additionally, comparing the performance of gene finders on other fungal genomes could provide insights into whether the findings are generalizable. Examining the performance of gene finders on different assemblies of the same genome could also be useful, as different assemblies may have different levels of completeness and accuracy. This could help to identify tangible benefits of using one assembly over another, and could also provide insights into the limitations of gene finders when applied to different assemblies.

Finally, the results of this work could be used to inform future research in the field of fungal genomics, particularly in the case of the DC1 and Tsth20 genome assemblies. The results of these genome assemblies are high-quality, and contain near-chromosomal scale sequences. The predicted genes from these genomes are a rich resource for future research, and could be used to identify novel genes and gene families in DC1 and Tsth20. Further understanding of the mechanisms behind salt and drought tolerance as well as degradation of hydrocarbon in suboptimal environments could be achieved by studying the predicted genes in these genomes.

References

- [1] Sheridan L. Woo, Rosa Hermosa, Matteo Lorito, and Enrique Monte. *Trichoderma*: A multipurpose, plant-beneficial microorganism for eco-sustainable agriculture. *Nature Reviews Microbiology*, 21(5):312–326, 2023.
- [2] Christian P. Kubicek, Andrei S. Steindorff, Komal Chenthamara, Gelsomina Manganiello, Bernard Henrissat, Jian Zhang, Feng Cai, Alexey G. Kopchinskiy, Eva M. Kubicek, Alan Kuo, Riccardo Baroncelli, Sabrina Sarrocco, Eliane Ferreira Noronha, Giovanni Vannacci, Qirong Shen, Igor V. Grigoriev, and Irina S. Druzhinina. Evolution and comparative genomics of the most common *Trichoderma* species. *BMC Genomics*, 20(1):485, 2019.
- [3] Prasun K. Mukherjee, Benjamin A. Horwitz, and Charles M. Kenerley. Secondary metabolism in *Trichoderma* – a genomic perspective. *Microbiology*, 158(1):35–45, 2012.
- [4] Carla Gonçalves, Chris Todd Hittinger, and Antonis Rokas. Horizontal Gene Transfer in Fungi and Its Ecological Importance. In Yen-Ping Hsueh and Meredith Blackwell, editors, *Fungal Associations*, pages 59–81. Springer International Publishing, 2024.
- [5] Miroslav Plohl, Nevenka Meštrović, and Brankica Mravinac. Centromere identity from the DNA point of view. *Chromosoma*, 123(4):313–325, August 2014.
- [6] David J. Winter, Austen R. D. Ganley, Carolyn A. Young, Ivan Liachko, Christopher L. Schardl, Pierre-Yves Dupont, Daniel Berry, Arvina Ram, Barry Scott, and Murray P. Cox. Repeat elements organise 3D genome structure and mediate transcription in the filamentous fungus *Epichloë festucae*. *PLOS Genetics*, 14(10):1–29, 2018.
- [7] Christian P. Kubicek, Alfredo Herrera-Estrella, Verena Seidl-Seiboth, Diego A. Martinez, Irina S. Druzhinina, Michael Thon, Susanne Zeilinger, Sergio Casas-Flores, Benjamin A. Horwitz, Prasun K. Mukherjee, Mala Mukherjee, László Kredics, Luis D. Alcaraz, Andrea Aerts, Zsuzsanna Antal, Lea Atanasova, Mayte G. Cervantes-Badillo, Jean Challacombe, Olga Chertkov, Kevin McCluskey, Fanny Culpier, Nandan Deshpande, Hans von Döhren, Daniel J. Ebbole, Edgardo U. Esquivel-Naranjo, Erzsébet Fekete, Michel Flippi, Fabian Glaser, Elida Y. Gómez-Rodríguez, Sabine Gruber, Cliff Han, Bernard Henrissat, Rosa Hermosa, Miguel Hernández-Oñate, Levente Karaffa, Idit Kost, Stéphane Le Crom, Erika Lindquist, Susan Lucas, Mette Lübeck, Peter S. Lübeck, Antoine Margeot, Benjamin Metz, Monica Misra, Helena Nevalainen, Markus Omann, Nicole Packer, Giancarlo Perrone, Edith E. Uresti-Rivera, Asaf Salamov, Monika Schmoll, Bernhard Seiboth, Harris Shapiro, Serenella Sukno, Juan Antonio Tamayo-Ramos, Doris Tisch, Aric Wiest, Heather H. Wilkinson, Michael Zhang, Pedro M. Coutinho, Charles M. Kenerley, Enrique Monte, Scott E. Baker, and Igor V. Grigoriev. Comparative genome sequence analysis underscores mycoparasitism as the ancestral life style of *Trichoderma*. *Genome Biology*, 12(4):R40, April 2011.
- [8] David A. Fitzpatrick. Horizontal gene transfer in fungi. *FEMS Microbiology Letters*, 329(1):1–8, 2012.
- [9] Biancamaria Senizza, Fabrizio Araniti, Simon Lewin, Sonja Wende, Steffen Kolb, and Luigi Lucini. *Trichoderma* spp.-mediated mitigation of heat, drought, and their combination on the *Arabidopsis thaliana* holobiont: A metabolomics and metabarcoding approach. *Frontiers in Plant Science*, 14, 2023.
- [10] Xiao-Ya An, Guo-Hui Cheng, Han-Xing Gao, Xue-Fei Li, Yang Yang, Dan Li, and Yu Li. Phylogenetic Analysis of *Trichoderma* Species Associated with Green Mold Disease on Mushrooms and Two New Pathogens on *Ganoderma Sichuanense*. *Journal of Fungi*, 8(7):704, July 2022.

- [11] N Saitou and M Nei. The neighbor-joining method: A new method for reconstructing phylogenetic trees. *Molecular Biology and Evolution*, 4(4):406–425, July 1987.
- [12] Joseph Felsenstein. Confidence Limits On Phylogenies: An Approach Using The Bootstrap. *Evolution; International Journal of Organic Evolution*, 39(4):783–791, July 1985.
- [13] Brendan Loftus. 4 - Genome Sequencing, Assembly and Gene Prediction in Fungi. In Dilip K. Arora and George G. Khachatourians, editors, *Fungal Genomics*, volume 3, pages 65–81. Elsevier, 2003.
- [14] Girum Fitihamlak Ejigu and Jaehee Jung. Review on the Computational Genome Annotation of Sequences Obtained by Next-Generation Sequencing. *Biology*, 9(9):295, September 2020.
- [15] Mark Borodovsky and Alex Lomsadze. Eukaryotic gene prediction using GeneMark.hmm-E and GeneMark-ES. *Current Protocols in Bioinformatics*, Chapter 4(SUPPL.35):4.6.1–4.6.10, September 2011.
- [16] C. Burge and S. Karlin. Prediction of complete gene structures in human genomic DNA. *Journal of Molecular Biology*, 268(1):78–94, April 1997.
- [17] W. H. Majoros, M. Pertea, and S. L. Salzberg. TigrScan and GlimmerHMM: Two open source *ab initio* eukaryotic gene-finders. *Bioinformatics*, 20(16):2878–2879, November 2004.
- [18] Mario Stanke and Stephan Waack. Gene prediction with a hidden Markov model and a new intron submodel. *Bioinformatics*, 19 Suppl 2:ii215–225, October 2003.
- [19] Tomáš Brůna, Katharina J. Hoff, Alexandre Lomsadze, Mario Stanke, and Mark Borodovsky. BRAKER2: Automatic eukaryotic genome annotation with GeneMark-EP+ and AUGUSTUS supported by a protein database. *NAR Genomics and Bioinformatics*, 3(1):lqaa108, March 2021.
- [20] Brandi L. Cantarel, Ian Korf, Sofia M.C. Robb, Genis Parra, Eric Ross, Barry Moore, Carson Holt, Alejandro Sánchez Alvarado, and Mark Yandell. MAKER: An easy-to-use annotation pipeline designed for emerging model organism genomes. *Genome Research*, 18(1):188–196, January 2008.
- [21] The NCBI Eukaryotic Genome Annotation Pipeline. https://www.ncbi.nlm.nih.gov/refseq/annotation_euk/process/, August 2025.
- [22] Mosè Manni, Matthew R. Berkeley, Mathieu Seppey, and Evgeny M. Zdobnov. BUSCO: Assessing Genomic Data Quality and Beyond. *Current Protocols*, 1(12):e323, 2021.
- [23] Arryn Craney, Salman Ahmed, and Justin Nodwell. Towards a new science of secondary metabolism. *The Journal of Antibiotics*, 66(7):387–400, 2013.
- [24] Hisayuki Komaki and Tomohiko Tamura. Polyketide Synthase and Nonribosomal Peptide Synthetase Gene Clusters in Type Strains of the Genus Phytohabitans. *Life*, 10(11):257, October 2020.
- [25] Kai Blin, Simon Shaw, Hannah E Augustijn, Zachary L Reitz, Friederike Biermann, Mohammad Alanjary, Artem Fetter, Barbara R Terlouw, William W Metcalf, Eric J N Helfrich, Gilles P van Wezel, Marnix H Medema, and Tilman Weber. antiSMASH 7.0: New and improved predictions for detection, regulation, chemical structures and visualisation. *Nucleic Acids Research*, 51(W1):W46–W50, July 2023.
- [26] Byoungnam Min, Igor V Grigoriev, and In-Geol Choi. FunGAP: Fungal Genome Annotation Pipeline using evidence-based gene model evaluation. *Bioinformatics*, 33(18):2936–2937, September 2017.
- [27] Simon Bennett. Solexa Ltd. *Pharmacogenomics*, 5(4):433–438, June 2004.
- [28] Anthony M. Bolger, Marc Lohse, and Bjoern Usadel. Trimmomatic: A flexible trimmer for Illumina sequence data. *Bioinformatics*, 30(15):2114–2120, August 2014.
- [29] Yunhao Wang, Yue Zhao, Audrey Bollas, Yuru Wang, and Kin Fai Au. Nanopore sequencing technology, bioinformatics and applications. *Nature Biotechnology*, 39(11):1348–1365, 2021.

- [30] Jiang Hu, Zhuo Wang, Zongyi Sun, Benxia Hu, Adeola Oluwakemi Ayoola, Fan Liang, Jingjing Li, José R. Sandoval, David N. Cooper, Kai Ye, Jue Ruan, Chuan-Le Xiao, Depeng Wang, Dong-Dong Wu, and Sheng Wang. NextDenovo: An efficient error correction and accurate assembly tool for noisy long reads. *Genome Biology*, 25(1):107, 2024.
- [31] Jiang Hu, Junpeng Fan, Zongyi Sun, and Shanlin Liu. NextPolish: A fast and efficient genome polishing tool for long-read assembly. *Bioinformatics*, 36(7):2253–2255, April 2020.
- [32] Simon Andrews. Babraham Bioinformatics - FastQC A Quality Control tool for High Throughput Sequence Data. <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.
- [33] Alexey Gurevich, Vladislav Saveliev, Nikolay Vyahhi, and Glenn Tesler. QUASt: Quality assessment tool for genome assemblies. *Bioinformatics*, 29(8):1072–1075, April 2013.
- [34] P. Rice, I. Longden, and A. Bleasby. EMBOSS: The European Molecular Biology Open Software Suite. *Trends in Genetics*, 16(6):276–277, June 2000.
- [35] Vijai Kumar Gupta, Andrei Stecca Steindorff, Renato Graciano de Paula, Rafael Silva-Rocha, Astrid R. Mach-Aigner, Robert L. Mach, and Roberto N. Silva. The Post-genomic Era of *Trichoderma reesei*: What’s Next? *Trends in Biotechnology*, 34(12):970–982, December 2016.
- [36] Chengzhi Liang, Long Mao, Doreen Ware, and Lincoln Stein. Evidence-based gene predictions in plant genomes. *Genome Research*, 19(10):1912–1923, October 2009.
- [37] SRA-Toolkit. <https://hpc.nih.gov/apps/sratoolkit.html>, June 2025.
- [38] Kolmogorov–Smirnov Test. In *The Concise Encyclopedia of Statistics*, pages 283–287. Springer New York, 2008.
- [39] Mosè Manni, Matthew R. Berkeley, Mathieu Seppey, and Evgeny M. Zdobnov. BUSCO: Assessing Genomic Data Quality and Beyond. *Current Protocols*, 1(12):e323, 2021.
- [40] Felipe A. Simão, Robert M. Waterhouse, Panagiotis Ioannidis, Evgenia V. Kriventseva, and Evgeny M. Zdobnov. BUSCO: Assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*, 31(19):3210–3212, October 2015.
- [41] Philip Jones, David Binns, Hsin-Yu Chang, Matthew Fraser, Weizhong Li, Craig McAnulla, Hamish McWilliam, John Maslen, Alex Mitchell, Gift Nuka, Sebastien Pesseat, Antony F. Quinn, Amaia Sangrador-Vegas, Maxim Scheremetjew, Siew-Yit Yong, Rodrigo Lopez, and Sarah Hunter. InterProScan 5: Genome-scale protein function classification. *Bioinformatics*, 30(9):1236–1240, 2014.
- [42] E. Michael Gertz, Yi-Kuo Yu, Richa Agarwala, Alejandro A. Schäffer, and Stephen F. Altschul. Composition-based statistics and translated nucleotide searches: Improving the TBLASTN module of BLAST. *BMC Biology*, 4(1):41, December 2006.
- [43] Rina Dechter. Chapter 12 - Temporal Constraint Networks. In Rina Dechter, editor, *Constraint Processing*, pages 333–362. Morgan Kaufmann, 2003.
- [44] Python Software Foundation. Python. <https://www.python.org/psf-landing/>, October 2025.
- [45] Brad Chapman. BCBio Python Package. <https://github.com/chapmanb/bcbio/tree/master/gff>, August 2025.
- [46] James T. Robinson, Helga Thorvaldsdóttir, Wendy Winckler, Mitchell Guttman, Eric S. Lander, Gad Getz, and Jill P. Mesirov. Integrative genomics viewer. *Nature Biotechnology*, 29(1):24–26, 2011.
- [47] Alexey Gurevich, Vladislav Saveliev, Nikolay Vyahhi, and Glenn Tesler. QUASt: Quality assessment tool for genome assemblies. *Bioinformatics*, 29(8):1072–1075, April 2013.

- [48] Mario Stanke, Ana Tzvetkova, and Burkhard Morgenstern. AUGUSTUS at EGASP: Using EST, protein and genomic alignments for improved gene prediction in the human genome. *Genome Biology*, 7(1):S11, 2006.
- [49] Derek W. Barnett, Erik K. Garrison, Aaron R. Quinlan, Michael P. Strömberg, and Gabor T. Marth. BamTools: A C++ API and toolkit for analyzing and managing BAM files. *Bioinformatics (Oxford, England)*, 27(12):1691–1692, June 2011.
- [50] Aaron R. Quinlan and Ira M. Hall. BEDTools: A flexible suite of utilities for comparing genomic features. *Bioinformatics*, 26(6):841–842, 2010.
- [51] Mihaela Pertea, Geo M. Pertea, Corina M. Antonescu, Tsung-Cheng Chang, Joshua T. Mendell, and Steven L. Salzberg. StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nature Biotechnology*, 33(3):290–295, 2015.
- [52] Geo Pertea and Mihaela Pertea. GFF Utilities: GffRead and GffCompare. *F1000Research*, 9, 2020.
- [53] Bradley W. Wright, Mark P. Molloy, and Paul R. Jaschke. Overlapping genes in natural and engineered genomes. *Nature Reviews Genetics*, 23(3):154–168, 2022.
- [54] S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman. Basic local alignment search tool. *Journal of molecular biology*, 215(3):403–410, October 1990.