Accepted Manuscript

Genomic analysis of the aggressive tree pathogen Ceratocystis albifundus

Magriet A. van der Nest, Emma T. Steenkamp, Danielle Roodt, Nicole C. Soal, Marike Palmer, Wai-Yin Chan, P. Markus Wilken, Tuan A. Duong, Kershney Naidoo, Quentin C. Santana, Conrad Trollip, Lieschen De Vos, Stephanie van Wyk, Alistair R. McTaggart, Michael J. Wingfield, Brenda D. Wingfield

PII: S1878-6146(19)30021-2

DOI: https://doi.org/10.1016/j.funbio.2019.02.002

Reference: FUNBIO 1002

To appear in: Fungal Biology

Received Date: 4 June 2018

Revised Date: 8 February 2019

Accepted Date: 11 February 2019

Please cite this article as: van der Nest, M.A., Steenkamp, E.T., Roodt, D., Soal, N.C., Palmer, M., Chan, W.-Y., Wilken, P.M., Duong, T.A., Naidoo, K., Santana, Q.C., Trollip, C., De Vos, L., van Wyk, S., McTaggart, A.R., Wingfield, M.J., Wingfield, B.D., Genomic analysis of the aggressive tree pathogen *Ceratocystis albifundus, Fungal Biology*, https://doi.org/10.1016/j.funbio.2019.02.002.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



1	Genomic analysis of the aggressive tree pathogen Ceratocystis albifundus
2	
3	Magriet A. van der Nest [†] , Emma T. Steenkamp, Danielle Roodt, Nicole C. Soal, Marike
4	Palmer, Wai-Yin Chan, P. Markus Wilken, Tuan A. Duong, Kershney Naidoo, Quentin C.
5	Santana, Conrad Trollip, Lieschen De Vos, Stephanie van Wyk, Alistair R. McTaggart,
6	Michael J. Wingfield, Brenda D. Wingfield
7	
8	Department of Biochemistry, Genetics and Microbiology, Forestry and Agricultural
9	Biotechnology Institute (FABI), University of Pretoria, Pretoria, South Africa.
10	
11	[†] Corresponding author:
12	Magriet A. van der Nest
13	E-mail: magriet.vandernest@fabi.up.ac.za
14	
15	Highlights
16	
17	• <i>Ceratocystis albifundus</i> has a two-speed genome.
18	• Its genome consists of core and accessory subgenomic compartments.
19	• Genetic variation is linked to the presence and activity of transposable elements.
20	• Genome structure likely contributed to pathogenicity and host specialization.
21	
22	
23	
24	

1 Abstract

2 Comparative genomics provides a powerful tool to investigate processes that underlie the 3 biology of fungi over evolutionary time. Such studies revealed that many pathogens have 4 "two-speed genomes", comprising of fast- and slow-evolving sub-genomic compartments. 5 The overall goal of this study was to determine whether the genome of an important plant 6 pathogen in Africa, Ceratocystis albifundus, is structured into sub-genomic compartments, 7 and if so, to establish how these compartments are distributed across the genome. For this 8 purpose, the publicly available genome of C. albifundus was complemented with the genome sequences for four additional isolates using the Illumina HiSeq platform. In addition, a 9 10 reference genome for one of the individuals was assembled using both PacBio and Illumina HiSeq technologies. Our results showed a high degree of synteny between the five genomes, 11 12 however several regions lacked detectable long-range synteny. These regions were associated with the presence of accessory genes, lower genetic similarity, variation in read-map depth, 13 as well as genes associated with host-pathogen interactions (e.g. effectors and CAZymes). 14 15 The regions lacking detectable long-range synteny were also abundant in transposable elements. The presence and activity of these elements is commonly associated with 16 accelerated evolution of accessory subgenomic compartments of fungal pathogens. Our 17 18 findings thus showed that the genome of C. albifundus is made-up of core and accessory sub-19 genomic compartments, which represents an important step towards characterizing its 20 pangenome (i.e., the combined core and accessory genomes). This study also highlights the value of comparative genomics as a tool to increase our understanding of the underlying 21 22 molecular mechanisms that may influence the biology and evolution of important pathogens. 23

Keywords: Ceratocystis wilt, Transposable elements, Pathogenomics, Adaptation, Host
jumps, Core and accessory genomes

- 26
- 27

1. Introduction

2 Genome comparisons provide a wealth of knowledge on the relationship among organisms 3 and the mechanisms that shape their biology (Hardison 2003; Wittenberg et al. 2009; 4 Goodwin et al. 2011). From a fungal perspective, genome comparisons have provided insights into the molecular basis of speciation, host-specificity and pathogenicity 5 mechanisms, as well as lineage-specific innovations (Plissonneau et al. 2017; Steenkamp et 6 7 al. 2018). For example, comparative genomic studies have revealed various genomic features 8 that are directly or indirectly responsible for species-specific lifestyle traits. These include 9 genome size and predicted gene products, the pathways and processes encoded by the genome, genome architecture, gain/loss of dispensable chromosomes, repetitive genomic 10 islands and genomic plasticity (Croll and McDonald 2012; Ma et al. 2013; Grandaubert et al. 11 12 2014; Shi et al. 2018).

13

1

Comparative genomics has revealed that many fungi have "two-speed genomes" made-up of 14 fast and slow-evolving sub-genomic compartments (Croll and McDonald 2012; Dong et al. 15 16 2015; Raffaele and Kamoun 2012). The slow-evolving compartment typically contains core 17 genes (i.e., those that are shared by all members of a species), the products of which mediate 18 general physiology and housekeeping functions. It is usually not very rich in repetitive 19 elements. This is in contrast to the fast-evolving compartment, which is typically repeat-rich, 20 architecturally dynamic and contains accessory genes (i.e., those that are absent from certain members of a species). In pathogens, accessory genes are often enriched for those involved in 21 22 virulence and host interactions (Ma et al. 2013; Faino et al. 2016; Plissonneau et al. 2016, 2018). Structurally, the fast-evolving subgenomic compartment may be distributed across all 23 24 or most chromosomes of an individual and/or reside on specific chromosomes that may be conditionally dispensable (Ma et al. 2010, 2013; Goodwin et al. 2011; Leclair et al. 1996; 25 Tzeng et al. 1992; Hatta et al. 2002; Plissonneau et al. 2018). 26

27

Previous genomic comparisons of diverse fungi have shown that transposable elements (TEs) play important roles in their evolution and adaptation (Grandaubert et al. 2014; Gladieux et al. 2014; Chiapello et al. 2015). TEs are mobile DNA segments capable of movement ("jumping") within a specific genome (Wicker et al. 2007; Amselem et al. 2015). This allows for the insertion of novel sequences within or close to existing genes, which may result in gene duplications, gene loss or gene inactivation (Daboussi 1996; Amselem et al. 2015; Biémont 2010). TEs are most prevalent in the fast-evolving subgenomic compartment (Dong

et al. 2015), and their activity affects the size, structure and dynamics of the genomes
harbouring them (Kidwell and Lisch 2000; Böhne et al. 2008). TE activity has been linked,
for instance, to accelerated evolution of genes involved in fungal pathogenicity and hostspecificity (Fudal et al. 2009; Manning et al. 2013). This is often due to the development of
gene repertoires implicated in niche expansion (Casacuberta and Santiago 2003) or whole
chromosomes enriched for TEs and genes associated with pathogenicity and virulence (Ma et
al. 2010; Goodwin et al. 2011).

8

In this study, we investigated the genomic substructure of *Ceratocystis albifundus* (Phylum: 9 Ascomycota; Order: Microascales; Family: Ceratocystidaceae) (De Beer et al. 2014). This 10 fungus is an aggressive pathogen of exotic Acacia mearnsii (Roux et al. 1999, 2001, 2005; 11 12 Heath et al. 2009) and commercially propagated Protea cynaroides in South Africa (Lee et al. 2016; Aylward et al. 2017). It has also been isolated from a wide range of native tree 13 species without causing obvious signs of disease (Roux et al. 2007). Even though C. 14 albifundus is homothallic with much of its reproduction occurring through selfing, 15 populations of this fungus have high levels of genetic diversity with intermediate levels of 16 gene flow (Roux et al. 2001, 2007; Barnes et al. 2005; Lee et al. 2016). This high diversity, 17 18 together with its wide host range and absence from other continents, suggests that C. albifundus is native in southern Africa (Roux et al. 2001, 2007; Barnes et al. 2005). Its 19 20 pathogenic niche on cultivated tree crops in South Africa may have resulted from a recent host jump and subsequent invasion (Roux et al. 2007). 21

22

Despite the importance of *C. albifundus* very little is known about the mechanisms 23 24 underlying its behaviour as an aggressive tree pathogen. Knowledge regarding its genomic make-up and substructure would inform our understanding of the molecular basis of its 25 26 biology, diversity and evolution. The overall goal of this study was, therefore, to determine 27 whether the genome of C. albifundus is structured into core and accessory compartments, and 28 if so, to establish how these compartments are distributed across the genome. We compared genes (especially those commonly associated with pathogenicity and interactions with the 29 30 plant host), TEs and repetitive elements, and analysed synteny across five genomes of C. albifundus from different hosts and geographic locations. For this purpose, the publicly 31 32 available genome of C. albifundus (van der Nest et al. 2014a) was complemented with 33 sequenced genomes of four additional isolates from a wide geographic range and different

hosts using Illumina HiSeq. In addition, a high-quality reference genome for one of the
 individuals was assembled using a combination of PacBio and Illumina HiSeq data.

3

2. Materials and Methods

4 5

2.1. Genome assemblies and annotation

Five isolates of C. albifundus originating from geographically diverse regions were included 6 7 in this study (Table 1). Three of the isolates were collected in South Africa (CMW 4068, 8 CMW 17274, CMW 17620), one in Zambia (CMW 13980) and one in Kenya (CMW 24685). Three of the isolates were from native trees (i.e., CMW 13980, CMW 17274 and CMW 9 17620) and two were from A. mearnsii (i.e., CMW 4068 and CMW 24685). Cultures were 10 obtained from the CMW Culture Collection at the Forestry and Agricultural Biotechnology 11 12 Institute at the University of Pretoria and maintained on 2 % malt extract agar (MEA, 20 gL-1 Agar, 20 gL-1 malt extract) at 25 °C (Lee et al. 2015). 13

14

The genome of isolate CMW 17620 was previously sequenced using the Illumina[®] platform 15 16 (van der Nest et al. 2014a). A total of 5 µg of DNA was prepared for the remaining isolates (CMW 4068, CMW 17274, CMW 13980 and CMW 17274) using previously described 17 methods (Barnes et al. 2001) and sequenced using the Illumina Genome Analyzer IIx 18 19 platform (Genome Centre, University of California, Davis, California, USA). CLC Genomics 20 Workbench v. 6.0.1 (CLC bio, Aarhus, Denmark) was used to trim the Illumina reads of low 21 quality (P error limit of 0.05). The remaining reads were assembled using the Velvet de novo 22 assembler with an optimal k-mer as determined with VelvetOptimiser (Zerbino 2010; 23 Zerbino and Birney 2008). Thereafter, the pre-assemblies were scaffolded using SSPACE v. 2.0 and gaps reduced using GapFiller v. 2.2.1 (Boetzer et al. 2011; Boetzer and Pirovano 24 2012). After assembly, the completeness of each genome was evaluated through the 25 26 Benchmarking Universal Single-Copy Orthologs (BUSCO) tool (BUSCO v. 1.1b1) by 27 determining the percentage of the most highly conserved fungal gene orthologs present in the 28 respective genomes (Simão et al. 2015). The genes or open reading frames (ORFs) for each 29 assembly were predicted using the *de novo* prediction software AUGUSTUS with Fusarium 30 graminearum gene models (Stanke et al. 2004) and then annotated using Blast2GO (Conesa 31 et al. 2005).

32

33 For isolate CMW 4068, we also assembled a high-quality hybrid reference genome using the

PacBio[®] RS II and Illumina HiSeqTM sequencing platforms. A total of 30 µg of genomic 1 2 DNA was prepared using the Qiagen DNeasy Plant Mini Kit (Qiagen, Valencia, CA, USA) 3 and sequenced by Macrogen (Seoul, Korea) using PacBio's Single Molecule Real Time (SMRT) sequencing technology. PacBio's SMRT Portal (v. 2.0.0) was used for read 4 corrections with the original PacBio parameters, after which the genome was assembled *de* 5 novo with Canu v. 1.7 (Berlin et al. 2015). For this, a range of values for the master errorRate 6 7 parameter was evaluated and a final errorRate of 0.0075 was used. Scaffolding was 8 performed using SSPACE-LongRead v. 1.1 (Boetzer and Pirovano 2014). By including the 9 quality-filtered Illumina data for isolate CMW 4068, the assembly was polished using Pilon v. 2 (Walker et al. 2014). For this hybrid assembly, completeness was estimated and ORFs 10 were predicted as before. 11

12

13

2.2. Identification and analysis of core and accessory compartments

To determine whether the genome of C. albifundus is separated into core and accessory 14 subgenomic compartments, nucleotide sequence and gene content was compared. For the 15 16 gene content-based analysis, we used the reciprocal protein Basic Local Alignment Search Tool (BLASTp) to determine whether individual genes were included in all or only some of 17 the Illumina genome assemblies. For this purpose, we used all the genes predicted from the 18 19 five Illumina assemblies. Those occurring in all five of the genomes were designated as "core genes", while those predicted to be present only in some of the isolates were designated as 20 21 "accessory genes". The latter also included the "unique genes" that were identified in only one of the isolates examined. All reciprocal BLASTp searches were done using a custom 22 23 python script (available from the authors). The script considered gene sequences with a 24 BLAST result Expect (E)-value cut-off of ≤ 0.00001 as shared between the pairs of genomes. 25 To control for annotation inconsistencies, individual translated nucleotide BLAST 26 (tBLASTn) searches (E-value cut-off of ≤ 0.00001) against the genomes were used to verify 27 that genes designated as accessory were indeed absent in one or more of the five assemblies 28 using BioEdit v. 7.2.5 (Hall 2011).

29

JSpecies (Richter and Rosselló-Móra 2009; Goris et al. 2007) was used to compare genomewide (coding plus non-coding) nucleotide similarity between pairs of the five Illumina
genome assemblies. For each pairwise analysis, JSpecies artificially sectioned one of the
genomes into fragments ranging between 100 and 1020 nucleotides in length, which were

then compared using the BLAST algorithm against the other genome, and vice versa (Altschul et al. 1997). In these comparisons, only those fragments that aligned over more than 70 % of their entire lengths and shared more than 30 % identity were considered as homologous (Goris et al. 2007). In addition to calculating the conserved fraction among the genomes, this software was also used to determine the average sequence similarity between genome pairs (Richter and Rosselló-Móra 2009; Goris et al. 2007).

7 8

2.3. Predicted pathways and processes for the core and accessory genes

9 The predicted proteins in the core and accessory datasets were mapped to the Kyoto Encyclopedia of Genes and Genomes (KEGG) databases using the GhostKoala mapping tool 10 11 (Kanehisa et al. 2016). GhostKoala assigned KEGG ORTHOLOGY identifiers (K numbers) 12 to each gene, which were used to reconstruct pathways on the KEGG web server (http://www.genome.jp/kegg/). For the accessory gene set, ClustVis (Metsalu and Vilo 2015) 13 14 was used to construct a copy number-based heatmap for each of the KEGG ORTHOLOGY identifiers in the five examined C. albifundus genomes. A Fisher's exact test (two-sided), 15 16 implemented in Blast2GO (Conesa et al. 2005), was employed to detect Gene Ontology (GO) 17 terms that were significantly enriched (P < 0.05) in the accessory set using the whole genome 18 as reference. The REVIGO web server (Supek et al. 2011) was used to summarize these 19 Blast2GO results.

20

The predicted proteins included in the core and accessory sets were also examined for the 21 22 presence of those known to be involved in interactions with their host, as well as virulence and pathogenicity related proteins (Ghannoum 2000; Paris et al. 2003; Tanabe et al. 2011; 23 24 Ohm et al. 2012; Zerillo et al. 2013; Li et al. 2015). These included Carbohydrate-Active 25 EnZymes (CAZymes), peptidases, catalases, superoxide dismutases, phospholipases and 26 effectors. Within the two datasets, we also identified genes whose products are potentially involved in the movement and activity of TEs. Chi-squared tests were used to determine 27 whether the frequency of these genes differed significantly (P < 0.05) between the core and 28 29 accessory gene sets (the null hypothesis was that the frequencies did not differ between the 30 two sets).

31

Putative CAZymes were identified and annotated with HMMER v. 3.0 (hmmer.org; Finn et
al. 2011) using the family-specific hidden Markov model profiles in the dbCAN database
(DataBase for automated Carbohydrate-active enzyme Annotation; Yin et al. 2012). Putative

1 peptidases were identified using the MEROPS database and BLASTp searches (E-value cutoff of ≤ 0.00001) (http://merops.sanger.ac.uk; Rawlings et al. 2016). A similar approach was 2 3 also used to compare genes across the five genomes to those in the Pathogen-Host Interactions database (PHI-base) v. 4.2 (http://www-phi4.phibase.org/), which includes a 4 5 collection of experimentally verified pathogenicity, virulence and effector genes from fungi, 6 oomycetes and bacteria (Winnenburg et al. 2008). To identify putative effectors, the two 7 datasets were first filtered for putative secreted proteins using SignalP v. 4.1 (Petersen et al. 8 2011), from which we then identified those with lengths, net charge and amino acid content typical of fungal effectors using EffectorP v. 1.0 (Sperschneider et al. 2016). Putative 9 proteins involved in the movement and activity of TEs were identified by BLASTp searches 10 (E-value cut-off of ≤ 0.00001) against 2636 known reference sequences from the "Core" set 11 of the Gypsy DataBase (GyDB) v. 2.0 (Llorens et al. 2011). 12

13

In each of the five genomes, putative catalases, superoxide dismutases and phospholipases 14 were identified using BLASTp searches (E-value cut-off of ≤ 0.00001) with previously 15 characterized protein sequences. These included catalases CATA, CATB, KATG1 and 16 KATG2 (NCBI accession numbers MG100061, MG06442, A4R5S9 and A4QUT2, 17 respectively), superoxide dismutases SOD1-5 (NCBI accession numbers EFZ03762, 18 EFY99820, EFY99375, EFZ00595 and EFZ00365, respectively) and phospholipases PLA2, 19 PLB2 and PLC2 (NCBI accession numbers KEY75421, KEY77760 and XP_011319931, 20 respectively). These sequences were aligned using MAFFT v. 7 (Multiple Alignment using 21 22 Fast Fourier Transform; http://mafft.cbrc.jp/alignment/server/) and then subjected to phylogenetic analysis using the distance-based neighbor-joining method in MEGA v. 7 23 24 (Molecular Evolutionary Genetic Analysis; http://www.megasoftware.net).

25

26 All putative host-associated genes included in the accessory gene set were subjected to 27 selection analysis with CODEML, as implemented in PAML v. 4.9h (Phylogenetic Analysis 28 by Maximum Likelihood; Yang and Nielson 2002). For this purpose, gene sequences from isolate CMW 4068 were used in local tBLASTx searches (i.e., translated nucleotide BLAST 29 30 searches against a translated nucleotide database) to identify homologs in the other four isolates. The identified sequences were then aligned with MAFFT and the phylogenetic trees 31 32 required by CODEML were inferred with MEGA as described above. Positive selection was 33 evaluated in each dataset by calculating the ratio (w) of non-synonymous (dN) versus 34 synonymous (dS) substitution rates across all sites (Yang et al. 2000). To test for variation of

selective pressures across the codons, goodness of fit was calculated for the different site specific models using likelihood ratio tests (Yang et al. 2000).

3

4

2.4. Genome conservation and syntemy

5 We compared the Illumina assemblies for isolates CMW 17620, CMW 17274, CMW 13980 6 and CMW 17274 against the PacBio-Illumina hybrid assembly for CMW 4068 using the 7 LASTZ (Large-Scale Genome Alignment Tool) (Harris 2007) and Mauve (Multiple 8 Alignment of Conserved Genomic Sequence With Rearrangements) (Darling et al. 2004) plugins implemented in Geneious v. 7 (Kearse et al. 2012). LASTZ aligned the Illumina 9 assemblies to the sequences of the 10 largest contigs in the hybrid assembly (contigs >1.1 10 Mb), by making use of "seed-and-extend" and "iterative refinement" strategies to allow 11 12 alignment of both the conserved and more variable regions (Harris 2007). LASTZ was also used for calculating the similarity between sequences. Mauve was used to plot sequence 13 14 similarity and to identify regions of local collinearity (Darling et al. 2004. The latter are known as Locally Collinear Blocks (LCBs), which are defined as homologous sequence 15 16 regions shared by the reference (i.e., the CMW 4068 hybrid assembly) and query (i.e., one of the Illumina assemblies), and that lack rearrangements. Synteny break points between the 17 18 reference and query genomes were further determined using SynChro (Drillon et al. 2014), 19 which employs Reciprocal Best-Hits (RBH) between coding sequences to identify conserved 20 syntenic blocks. The pairwise alignments extracted from LASTZ were annotated with AUGUSTUS and used as input for the SynChro analysis. SynChro then computed RBH to 21 22 reconstruct synteny block backbones, after which it automatically completed these blocks with non-RBH syntenic homologs. 23

24

CLC Genomics Workbench was used to determine the genomic distribution of single 25 26 nucleotide polymorphisms (SNPs). The quality-filtered Illumina reads for each isolate were 27 mapped (using default parameters) to the sequences for the 10 largest contigs in the CMW 28 4068 hybrid assembly. The R/Bioconductor package KaryoplotEr (Gel and Serra 2017) was used to calculate and plot SNP distribution in 5000 bp, non-overlapping intervals across the 29 30 sequence by expressing SNP content as the number of SNPs per interval. This package was also used to examine coverage or read depth using a sliding window of 5000 bp. The latter 31 32 was calculated following the removal of duplicate reads (to avoid over-representation of 33 specific reads due to sample preparation artefacts), with the maximum representation of a 34 minority sequence set to 20 %.

- 1
- 2 3

2.5. Genomic distribution of host-associated and accessory genes in the CMW 4068 hybrid assembly

4 Synteny and sequence similarity were used to identify genes in the CMW 4068 hybrid 5 assembly that potentially forms part of the accessory genome of C. albifundus. For this purpose, we aligned the individual Illumina assemblies to the hybrid genome assembly using 6 7 LASTZ (see above). In these alignments, genes that were missing in one or more of the 8 Illumina genome sequences were regarded as accessory genes in the CMW 4068 assembly. 9 The locations of these accessory genes were plotted across the 10 largest contigs in the hybrid 10 assembly using KaryoplotEr. Genes potentially involved in host interactions (CAZymes, peptidases and effectors) were identified as described before and their locations also plotted 11 12 with KaryoplotEr. Differences in distribution of these accessory genes among the different genomic regions were evaluated using Chi-squared tests as described above. 13

14

15 16

2.6. TE identification, annotation and distribution in the CMW 4068 hybrid assembly

17 The TEdenovo and TEannot pipelines in the REPET v. 2.3 package (Flutre et al. 2011) were 18 used to identify and annotate TEs in the CMW 4068 hybrid assembly. TEs were detected de *novo* with tBLASTx using E-value thresholds ($E < 10^{-10}$) (Gish and States 1993), the 19 BLASTER suite using similarity thresholds (E = 10^{-300} , minimum identity = 90 %) 20 (Quesneville et al. 2003), and LTRHarvest using TE structure (Ellinghaus et al. 2008). The 21 22 TEdenovo pipeline was then used to cluster the identified TEs, and to reconstruct a consensus for each group of matches using the programs Piler (Edgar and Myers 2003), GROUPER 23 24 (Quesneville et al. 2003) and RECON (Bao and Eddy 2002). Cut-offs for clustering individual TEs were set at 90 % identity over 95 % of the length (Flutre et al. 2011). 25 Consensus TE sequences were classified using the Repbase Update database 26 27 (http://www.girinst.org/repbase/update/index.html) and named according to the classification 28 proposed by Wicker et al. (2007).

29

30 Genomic locations identified TEs of the were plotted using KaryoploteR 31 (https://bioconductor.org/packages/release/bioc/html/karyoploteR.html) across the 10 largest 32 contigs of the hybrid assembly. This software was also used to plot the location of genes 33 potentially involved in the movement and activity of TEs in isolate CMW 4068, which were 34 identified using BLASTp searches against the "Cores" set of GyDB as described above. For

the latter, differences in distribution among specific genomic regions were evaluated with
 Chi-squared tests.

3

5

4 **3. Results**

3.1. Genome assemblies and annotation

After data filtering, we generated high-quality raw Illumina sequence data for C. albifundus 6 7 isolate CMW 4068 (mean read length of 94.6 bases), isolate CMW 13980 (mean read length 8 of 95.7 bases), isolate CMW 17274 (mean read length of 88.8 bases) and isolate CMW 24685 (mean read length of 95.7 bases) (Supplementary Table 1). The filtered data were used to 9 generate four genome assemblies that consisted of 818–2122 contigs and that were 26.7–27.2 10 Mb in size (Table 2), which is similar to the 27.3 Mb-assembly published for isolate CMW 11 12 17620 (van der Nest, et al. 2014a). Based on the BUSCO results (Simão et al. 2015), all five of the genomes were more than 98.0 % complete (Supplementary Table 2), which is 13 congruent with the general trend observed for *Ceratocystis* genomes (van der Nest et al. 14 15 2014a, 2014b, 2015; Wilken et al. 2013; Wingfield et al. 2015, 2016a, 2016b).

16

A total of 443 868 quality-filtered PacBio sequence reads (with an average length of 8 317 17 18 bases) was generated for CMW 4068 (Supplementary Table 1). Assembly, scaffolding and polishing yielded 16 contigs larger than 200 000 bases, spanning a total of 28.36 Mb and 19 20 containing 7 103 genes (Table 2 and Supplementary Table 3). Compared to the Illumina assemblies, the hybrid assembly had a much higher N50-value (2.31 Mb compared to 0.02-21 22 0.07 Mb) and size for its largest contig (4.1 Mb compared to 0.15-0.36 Mb). The hybrid PacBio-Illumina assembly and the Illumina assemblies were similar regarding gene density 23 24 (250 genes/Mb) and BUSCO completeness (98 %) (Table 2 and Supplementary Table 2), 25 which allows for meaningful gene-based genomic comparisons.

26

The CMW 4068 hybrid assembly was less fragmented than the Illumina assemblies (Table 2). This is because PacBio long-read sequencing allowed for the closing of gaps and sequencing through repetitive regions (Yue et al. 2017). However, the assembly included an additional 27 contigs, each of which consisted of less than 200 000 bases. In total, the 27 small contigs spanned 1.36 Mb and contained 341 genes, of which a significantly large portion (42 %; Chi-squared test, P > 0.05) showed similarity to genes in the GyDB database with a potential role in TE activity and movement (Llorens et al. 2011). Another 21 % shared

similarity to known genes, while the remaining genes (37 %) did not show similarity to
 known genes in any public database.

- 3
- 4

3.2. Identification and analysis of core and accessory compartments

5 The Illumina assemblies for the five isolates of C. albifundus were predicted to share 6 241 6 genes, which were included in the core set. This represented a large portion (92.4–95.4 %) of 7 the total number of genes predicted in each assembly (Figure 1, Supplementary Table 4). The 8 accessory genes (i.e., those missing from one or more of the genomes) numbered between 300 (representing 4.6 % of the CMW 17620 genome) and 515 genes (representing 7.6 % of 9 the CMW 24685 genome) (Supplementary Tables 5 and 6). A small proportion of accessory 10 genes were identified in only a single isolate (i.e., "unique genes"). These ranged from 17 11 12 (0.3 % of the genes encoded on the genomes of isolates CMW 4068 and CMW 24685) to 51 genes (0.8 % of genes encoded on the CMW 17274 genome) (Figure 1, Supplementary Table 13 14 6).

15

Genome-wide nucleotide comparisons of the 5 Illumina assemblies revealed that the 16 conserved fraction of the five C. albifundus genomes was very similar (Figure 2). Based on 17 18 the simple BLAST-based alignment strategy implemented in JSpecies, the average sequence 19 similarity values for the various pairwise comparisons ranged from 96.00 % (for isolates CMW 17620 and CMW 13980) to 98.63 % (for isolates CMW 24685 and CMW 4068). The 20 non-conserved or variable regions (i.e., where JSpecies fragments that did not align across 70 21 22 % or more of their lengths and that were not 30 % or more similar) represented between 3.31 % and 14.96 % of the total genomes of these fungi. These variable regions may represent 23 24 parts of the genome that were either too repetitive to be included in our assemblies with Velvet and/or that represent parts of the genome that are missing or variable among the 25 26 genomes compared.

27 28

3.3. Predicted pathways and processes for the core and accessory genes

Only a portion of core (46.1 %, Table 3, Supplementary Table 7) and accessory (12–21 %; Supplementary Table 8; Supplementary Figure S1) genes could be assigned to specific orthologs in the KEGG database. Among the various KEGG categories identified for the core genes, some are likely involved in housekeeping functions (e.g., "Transcription" and "Replication and repair"). However, it is possible that some categories are involved in functions related to niche utilization or host-pathogen interactions (e.g., "Biosynthesis of

other secondary metabolites" and "Metabolism of terpenoids and polyketides"). As was
expected, several of the accessory genes may be linked with host-pathogen interactions. Of
the KEGG categories identified for the accessory genes, some are likely related to niche
utilization or host-pathogen interactions (e.g., "Carbon metabolism", "Fructose and mannose
metabolism", "Fatty acid metabolism", "Histidine metabolism", "MAPK signalling pathway",
"mTOR signalling pathway", and "Ras signalling pathway") (Supplementary Table 8).

7

8 Fisher's exact test indicated that numerous GO terms were overrepresented in the accessory gene set (P < 0.05). These included GO terms associated with processes potentially involved 9 in niche utilization and/or host-pathogen interactions (Supplementary Table 9; 10 Supplementary Figure S5). The enriched GO terms were involved in biological functions 11 12 (e.g., "Establishment or maintenance of actin cytoskeleton polarity", "Establishment or maintenance of cell polarity", "Establishment or maintenance of cytoskeleton polarity") and 13 cellular components associated with host penetration (e.g., "New growing cell tip" and "Old 14 15 growing cell tip"). We also found enrichment in GO terms associated with secondary metabolism (e.g., "Aromatic compound biosynthetic process" and "Cellular aromatic 16 compound metabolic process"), cellular transport (e.g., "Nitrogen compound transport" and 17 18 "Import into cell" and signalling (e.g., "Cdc42 protein signal transduction").

19

Further detailed analyses showed that several of the core genes may be linked with host-20 pathogen interactions. These included genes (1095) that shared significant similarity with 21 22 previously characterized pathogenicity associated genes in the PHI-database (e.g., "fungal development, secondary metabolism and virulence" and "fungal development and 23 24 pathogenicity") (Supplementary Table 10), as well as genes (56) that encoded putative effectors (Supplementary Table 4) that may modulate the host immune system and promote 25 infection (Stergiopoulos and De Wit 2009). Among the core genes, the MEROPS-based 26 27 analyses identified various putative peptidases (260) involved in protein degradation and 28 modification (Supplementary Table 11). CAZymes (243) potentially involved in the degradation of plant polysaccharide materials (Supplementary Table 12) and phospholipases 29 30 potentially capable of hydrolysing plant phospholipids (Supplementary Figure S2) for facilitating infection and/or gaining nutrition were also identified among the core genes 31 32 (Cantarel et al. 2009; Ghannoum 2000; Zhao et al. 2013). The core genes further included 33 putative catalases (Supplementary Figure S3) and superoxide dismutases (Supplementary

1 Figure S4) that are known to scavenge reactive oxygen species to protect fungi from the host

- 2 defence responses (Tanabe et al. 2011; Li et al. 2015).
- 3

4 Many accessory genes (222) also shared significant similarity with previously characterized 5 pathogenicity associated genes in the PHI-database, including genes predicted to be involved 6 in "Cell wall adhesion", "Establishment of turgor in appressoria" and "Appressorial 7 penetration" (Supplementary Table 10). The accessory gene set also included those encoding 8 putative CAZymes (Supplementary Table 12). These included enzymes that catalyse the degradation of chitin (families CBM18 and GH18) (Hartl et al. 2012) and lignin (family 9 AA7) (Levasseur et al. 2013), as well as putative carbohydrate esterases responsible for 10 deacetylating plant polysaccharides (families CE4 and GH10) (Biely 2012). The accessory 11 12 gene set further contained genes encoding putative peptidases (Supplementary Table 11) with possible roles in fungus-plant interactions (e.g., two putative proteases in the sedolisin family 13 [Serine peptidase family \$53], and eleven putative subtilisin peptidases [Serine peptidase 14 family S08]) (Rawlings et al. 2016). Another 18 genes encoded putative M13 peptidases 15 (Metallo peptidase family M13) that likely play a role in regulation of peptide signalling 16 (Bland et al. 2008). Even though, the core and the accessory compartments of C. albifundus 17 encoded significantly different gene repertoires (Chi-squared test, P > 0.05), the frequency of 18 specific pathogenicity-related genes (i.e., peptidase, CAZyme and PHI-base genes) did not 19 20 differ significantly between the core and accessory gene sets (Chi-squared test, P > 0.05).

21

The accessory gene set included a large number of genes encoding putative effectors 22 (Supplementary Tables 5 and 6). These differed greatly among the isolates examined, ranging 23 from 11 putative effector genes in isolate CMW 17620 to 36 in isolate CMW 13890. The 24 frequency of these in the accessory gene set differed significantly from that in the core set 25 (Chi-squared test, P-value < 0.05). Furthermore, a number of these genes also appeared to 26 27 evolve under diversifying selection (Supplementary Table 16), because CODEML indicated 28 that the positive-selection models (M5, M6 and M8) provided a better fit compared to those 29 that assume no positive selection (M1 and M7) (Yang et al. 2000). This was also true for 30 other host-associated genes present in the accessory set, i.e., putative peptidases (M13 peptidases implicated in regulation of peptide signalling) and CAZymes (families CBM18 31 32 and GH18 involved in chitin degradation).

33

1 A large number of core and accessory genes were predicted to encode proteins associated with the activity of TEs. For the core gene set, searches against the reference sequences in 2 3 GyDB database, recovered 1009 putative genes involved in TE activity or their movement 4 (Supplementary Table 13). Likewise, the accessory gene set included 36-80 sequences with 5 similarity to GyDB database genes (e.g., genes encoding reverse transcriptases, retroelement integrases and *Gag*-like proteins involved in the replication and integration of certain TEs) 6 7 (Wilhelm and Wilhelm 2001; Novikova 2009; Llorens et al. 2011). However, the frequency 8 of these genes differed significantly between the core and accessory gene sets (Chi-squared test, P < 0.05) for isolates CMW 4068, CMW 24685 and CMW 13890. Also, Fisher's exact 9 test indicated that the accessory gene set was significantly (P < 0.05) enriched in functions 10 related to DNA integration processes that are known to play a role in the insertion of 11 transposable elements in protein coding genes (Plissonneau et al. 2018). 12

13

14

3.4. Genome conservation and synteny

Alignments of the Illumina assemblies against the 10 largest contigs of the hybrid assembly 15 for isolate CMW 4068 showed interrupted stretches of high similarity. Similarity between the 16 Illumina assemblies and the hybrid assembly ranged from 96.2 % for contig 7 of CMW 17 18 13890 to 99.4 % for contig 1 of CMW 17620 (Supplementary Table 14). On most of the contigs, large stretches of similarity (i.e., LCBs identified with Mauve) were interrupted by 19 20 areas that were unalignable. This was particularly evident on contigs 1, 4 and 7 (Figure 3 and Supplementary Figure S6). Also, a fraction of the individual Illumina reads for CMW 13890 21 22 (1.9 %), CMW 17620 (1.1 %), CMW 17274 (2.9 %) and CMW 24685 (1.5 %) did not map to the hybrid assembly for isolate CMW 4068 (Supplementary Table S15). The proportion of 23 24 these "unmapped" fractions was generally somewhat lower than those observed using pairwise BLAST comparisons of the Illumina assemblies (Fig. 2). This may be because the 25 repetitive nature of the individual genomes influences how many reads map to the reference 26 27 genome, although we also cannot exclude the possible influence of differences in genome 28 completeness and genome size on these data.

29

Overall, a high level of synteny (Figure 4A) was observed among the five *C. albifundus*genomes based on the presence and order of orthologous genes using SynChro (Drillon et al.
2014). Despite this high level of gene order conservation, numerous regions lacking
detectable long-range synteny were also observed relative to the 10 largest contigs of the
CMW 4068 hybrid assembly (Figures 3, 4 and Supplementary Figure 2). SynChro detected

synteny breaks (>11.1 kb on average) on contigs 1, 4, 6, 7, 8 and 9 (Figure 3 and
 Supplementary Figure S6). Additionally, SynChro detected inversions between the hybrid
 assembly and Illumina assemblies: three inversions for CMW 17274, 7 inversions each for
 CMW 17620 and CMW 13890, and 8 inversions for CMW 24685.

5

6 The regions lacking detectable long-range synteny often occurred in SNP-dense regions 7 and/or regions with large variation in read depth (Figure 3; Supplementary Figure S6). For 8 example, the three small synteny breaks (from positions 0.53 to 0.55 Mb, 0.63 to 0.67 Mb, and 1.92 to 1.96 Mb in the hybrid assembly) and one large synteny break (positions 1.11 to 9 1.36 Mb in the hybrid assembly) on contig 7 were all localized in areas characterized by a 10 higher number of SNPs in one or more of the compared genomes. The fact that these breaks 11 12 typically occurred in regions with highly variable read depth (e.g., the read depth for the large syntenic break on contig 7 ranged from 20 to 170 reads/5000 bases among the genomes 13 compared) suggests that the breaks were primarily caused by the presence of repetitive 14 15 sequences in these areas.

- 16
- 17 18

3.5. Genomic distribution of host-associated and accessory genes in the CMW 4068 hybrid assembly

Relative to the total number of genes encoded per contig, the abundance of accessory genes 19 20 (i.e., those that were missing in one or more of the four Illumina assemblies) varied substantially across the 10 largest contigs of the hybrid assembly (Supplementary Figure S6). 21 22 For example, accessory genes were much more abundant on contigs 6 (15.1 %) and 7 (19.2 %) than contigs 2 (8.9 %), 3 (5.4%) and 8 (8.9 %) (Figure 3 and Supplementary Figure S6). 23 24 Although the accessory genes were generally distributed across contigs, they often appeared to be localized in regions lacking detectable long-range synteny and that SynChro identified 25 as syntenic breaks (e.g., positions 0.78 to 0.89 Mb on contig 4, positions 0.13 to 0.50 Mb and 26 27 1.94 to 1.99 Mb on contig 6 and on contig 7 positions 1.92 to 1.96 Mb). Accordingly, Chi-28 squared tests of independence rejected the null expectation that the frequency of accessory 29 genes located in regions lacking long-range synteny is the same as in the rest of the contig (P-30 value < 0.05).

31

The genes encoding products potentially involved in host interactions and in the activity and movement of TEs appeared to be randomly positioned on the 10 largest contigs of the hybrid assembly. These included CAZymes, peptidases and putative effectors, as well as GyBD-

identified genes involved in movement and activity of TEs (Figure 3, Supplementary Figure
S6 and Supplementary Tables 11 - 13). However, the non-syntenic regions appeared to be
enriched for putative effectors (contig 7), CAZyme (contigs 6 and 8) and TE-associated
(contigs 1, 4 and 6) genes (Figure 3 and Supplementary Figure S6).

5

6

7

3.6. TE identification, annotation and distribution in the CMW 4068 hybrid assembly

Based on their predicted transposition mechanisms, the transposable elements and repeat 8 9 sequences for the hybrid assembly were classified as Class I, Class II or "NoClass" for those that could not be assigned into either of the two classes by REPET (Figure 4B; 10 Supplementary Table 17). Most of the annotated TEs represented Class I transposons 11 12 (commonly referred to as retrotransposons), which utilize a copy-and-paste mechanism for transposition via an RNA intermediate (Wicker et al. 2007; Amselem et al. 2015). Those TEs 13 14 annotated as being Class II transposons (commonly referred to as DNA transposons) likely 15 use a cut-and-paste mechanism involving transposases and DNA intermediates (Wicker et al. 2007; Amselem et al. 2015). The annotated Class I and Class II TEs were further grouped 16 into orders (Wicker et al. 2007), based on their insertion mechanism, structure and encoded 17 proteins (Figure 4C; Supplementary Table 17). The annotated TEs were represented by 5 18 orders, namely LTR [Long Terminal Repeat], LINE [Long Interspersed Nuclear Elements], 19 TIR [Terminal Inverted Repeats], DIRS-like elements [Dictyostelium intermediate repeat 20 sequence] and PLE [Penelope-like elements]. 21

22

A large portion (16.3 % for the 16 contigs > 2 Mb) of the hybrid assembly is represented by 23 24 TEs. The occurrence and distribution of the various TEs differed substantially among contigs (Figures 3, 4B, 4C, Supplementary Table 17). For instance, contig 7 contained a much higher 25 proportion of TEs (30 %) than contig 2 of which only 5.5 % were dedicated to TEs. The more 26 27 TE-dense regions also appeared to occur in areas lacking detectable long-range synteny 28 (Figure 3A and Supplementary Figure 2). On contig 7, for example, all four of the syntenic 29 breaks co-localize with TE-dense regions (Figure 3A). In terms of TE distribution, we 30 observed some clustering (e.g., in syntenic breaks and at the ends of the contigs), however, many TEs also seemed to be spread out across contigs (Figure 3A and Supplementary Figure 31 32 2).

33

1 The ends of almost all of the contigs in the CMW 4068 hybrid assembly were rich in TEs and repeats. Chi-squared tests showed that the frequencies of these elements within the terminal 2 3 20 000 bases of each contig were significantly different from those outside these regions (P <4 0.05). Furthermore, within the terminal 20 000 bases, five of the large contigs (i.e., 4, 5, 6, 12 5 and 14) contained 20-38 copies of the telomeric repeat 5' TTAGGG 3' motif (Fulnečková et 6 al. 2013; van Wyk et al. 2018), but only at one of their ends. It is therefore unlikely that any 7 of our contigs represent chromosome-sized scaffolds, indicating that the Pacbio long-reads 8 were not adequate for sequencing across the repetitive regions and assembling to 9 chromosome level.

10

11

4. Discussion

The results of this study demonstrated that the genome of C. albifundus is comprised of core 12 and accessory subgenomic compartments. This is similar to what has been observed in other 13 fungi (Croll and McDonald 2012; Gladieux et al. 2014; Ma et al. 2013; Dong et al. 2015; 14 Ohm et al. 2012 and may correspond to the eu- and heterochromatic DNAs of eukaryotes 15 (Vanrobays et al. 2018). In our study, this was evidenced by the stretches of individual 16 genomes that were highly variable and lacking synteny, and that were interspersed with 17 conserved regions of high sequence similarity. Comparisons of our five Illumina assemblies 18 19 also revealed large proportions of genes common to all isolates of C. albifundus (viz. forming part of the core subgenomic compartment), as well as genes (4.6–7.6 % of those predicted per 20 isolate) that were present in only some isolates (viz. forming part of the accessory 21 subgenomic compartment). Such individual or lineage-specific patterns of gene 22 presence/absence have also been reported for the heterochromatic regions of other eukaryotes 23 24 (Fortna et al. 2004; Dopman and Hartl 2007). Our study thus represents an important step towards characterizing the pangenome (i.e., the combined core and accessory genomes) of C. 25 26 albifundus, as these subgenomic compartments have important and distinctly different roles 27 in the biology and evolution of a pathogen (Plissonneau et al. 2018).

28

The core and the accessory compartments of *C. albifundus* encoded different gene repertoires, consistent with previous reports (Plissonneau et al. 2016, 2018). Core genes were predicted to mostly encode basal or housekeeping functions (see Table 2), while accessory genes encoded putative products required for access to plant-derived nutrients (Lee and Sheppard 2016; Ohm et al. 2012; Zhao et al. 2013) and host-pathogen interactions (Desjardins and Hohn 1997; Mukherjee et al. 2012). The accessory genes also encoded

putative products for signal transduction in sensing and responding to environmental conditions, thus enabling growth and survival in a specific biological niche (Braunsdorf et al. 2016). This knowledge provides a foundation for future functional studies that aim to clarify the roles of these proteins, and the possibility of manipulating them to improve disease management.

6

7 Structural analyses of the C. albifundus genomes suggested that the accessory subgenomic regions are distributed throughout the genome rather than located on specific chromosomes 8 (Ma et al. 2013; Goodwin et al. 2011; Leclair et al. 1996; Tzeng et al. 1992; Hatta et al. 9 2002). The percentage and distribution of accessory genes in the genome is comparable to 10 what has been reported in *Zymoseptori tritici* (Plissonneau et al. 2016, 2018). The accessory 11 subgenomic compartment of C. albifundus may contain 4.9-8.3 % of all the genes predicted 12 in an isolate (in Z. tritici 1.8-8.5 % of genes lack homologs in one or more isolates) 13 (Plissonneau et al. 2016, 2018). Like in Z. tritici, many C. albifundus accessory genes also 14 co-occurred in areas lacking long-range synteny (Plissonneau et al. 2016, 2018). However, 15 Ceratocystis differs from fungi such as Fusarium, where isolate or lineage-specific genomic 16 17 regions are localized to specific chromosomes (Ma et al. 2010, 2013).

18

The accessory subgenomic compartments in fungi are often enriched for genes involved in 19 virulence and pathogenicity (Ma et al. 2013; Faino et al. 2016; Plissonneau et al. 2016, 2018). 20 21 This may also be true for *C. albifundus*, since many non-syntenic regions contained genes 22 encoding putative CAZymes and effectors (Ohm et al. 2012; Dong et al. 2015; Fouché et al. 2018; Laurie et al. 2012; Faino et al. 2016). Effector genes commonly reside in rapidly 23 24 evolving accessory compartments in genomes of filamentous plant pathogens (Laurie et al. 2012; Dong et al. 2015; Faino et al. 2016; Fouché et al. 2018). As shown before, putative 25 26 effector genes were identified in the accessory compartment of C. albifundus. Also, 27 consistent with what has been observed for the accessory genes in other fungal pathogens 28 (Raffaele and Kamoun 2012), some of the host-associated functions (i.e., encoding effectors, CAZymes and peptidases) encoded on the accessory set of C. albifundus also evolve under 29 30 diversifying selection. These data suggest that some of the accessory genes in C. albifundus encode products that modulate host immune responses, promote infection and effective 31 32 colonization of plant hosts and that allow adaptation of the fungus to changing environments 33 (Stergiopoulos and de Wit 2009).

34

1 Our study suggests that C. albifundus has a two-speed-genome (Croll and McDonald 2012; Dong et al. 2015). This is because the regions enriched with accessory genes in the C. 2 3 *albifundus* genome appeared to be co-localised with synteny breaks, and these breaks bear all 4 of the hallmarks of fast-evolving regions. They contained unique sequences and genes, 5 displayed low sequence similarity and high SNP density, as well as vary greatly in Illumina sequencing read depth. This is similar to what has been observed in Verticillium dahliae 6 7 (Faino et al. 2016). Also, disruptions in long-range synteny have been linked to the divergence of Saccharomyces cerevisiae and its wild relative S. paradoxus (Yue et al. 2017). 8 In C. albifundus, the presence of such a fast-evolving subgenomic compartment may help to 9 10 explain the host range and high levels of genetic diversity reported in previous studies (Roux 11 et al. 2001, 2007; Barnes et al. 2005).

12

In fungal pathogens, the accelerated evolutionary rates of accessory subgenomic 13 14 compartments are often linked to the presence of TEs (Croll and McDonald 2012; Raffaele and Kamoun 2012; Vanheule et al. 2016; Faino et al. 2016). Consistent with this view, the 15 accessory genes of C. albifundus were overrepresented (38-65 %) for genes that encode 16 proteins involved in the movement or activity of TEs. Also, the genomic locations of synteny 17 18 breaks and TE density were clearly coordinated. The TE content of C. albifundus (16.3 % for 19 the hybrid reference genome) was higher than what was reported for many other fungi. These 20 include other necrotrophic pathogens such as Botrytis cinerea (0.7-2.2 %), Stagonospora nodorum (2.4 %), Alternaria brassicicola (5.6 %), as well as hemibiotrophs like Dothistroma 21 22 septosporum (0.7 %), Cochliobolus sativus (5.4 %) and Mycosphaerella populorum (3.6 %) (Grandaubert et al. 2014; Ohm et al. 2012; Amselem et al. 2011, 2015). In fact, the TE 23 24 content in the C. albifundus genome was similar to fungi in which large-scale TE invasions have been reported. These include hemibiotrophs such as *Leptosphaeria maculans* (25 %) 25 26 and *Mycosphaerella graminicola* (11.7 %), as well as ectomycorrhizal fungi such as *Laccaria* 27 bicolor (24 %) and Tuber melanosporum (58 %) (Ohm et al. 2012; Labbe et al. 2012; 28 Raffaele and Kamoun 2012; Grandaubert et al. 2014).

29

The findings presented in this study suggest that TEs likely played a significant role in the evolution of *C. albifundus*. The genome of this fungus is relatively rich in these elements and some of the genetic diversity observed in *C. albifundus* might have been the product of varied TE activities over time. Through their activity, TEs could have shaped the genomic landscape of *C. albifundus* by causing chromosomal rearrangements, deletions and duplications

(Daboussi 1997; Daboussi and Capy 2003). In addition to gene and genomic plasticity, TE
activity might also have caused phenotypic diversity through epigenetic mechanisms and/or
other changes in gene regulation (Daboussi and Capy 2003). TEs and the dynamic accessory
subgenomic compartment in which they occur are thus likely to have been important sources
of diversity and adaptive phenotypes in *C. albifundus* (Croll and McDonald 2012; Gladieux
et al. 2014; Ma et al. 2013; Dong et al. 2015).

7

8 Acknowledgements

We thank the University of Pretoria, the Department of Science and Technology (DST) 9 10 National Research Foundation (NRF) Centre of Excellence in Tree Health Biotechnology, the DST-NRF SARCHI Chair in Fungal Genomics, the Genomics Research Institute at the 11 12 University of Pretoria and the Claude Leon Foundation (South Africa) for financing the sequencing of the genomes included in this study. Financial support from The Improving 13 14 Academic Oualifications (IAO) NRF Sabbatical grant also contributed to cost of sequencing. This study was based on research supported by multiple grants from the NRF of South 15 16 Africa, and these include Grant specific unique reference number (UID) 83924. The Grant holders acknowledge that opinions, findings and conclusions or recommendations expressed 17 in publications generated by NRF supported research are that of the authors, and that the 18 19 NRF accepts no liability whatsoever in this regard.

20

21 Data Availability

The Whole Genome Shotgun projects have been deposited at DDBJ/EMBL/GenBank under accession numbers MAOA0000000, MANZ0000000, MANY00000000, MANX00000000 and MANW00000000. The hybrid assembly for isolate CMW 4068 was used to update the existing Illumina sequence for this isolate, and is available as version MAOA02000000.

1 Figure legends

2

Figure 1. The proportion of genes predicted to be encoded by the examined *Ceratocystis albifundus* genomes relative to the proportion of genes common to all five (i.e., the core genes), as well as the genes unique to a specific isolate and the genes absent from one or more of the isolates (i.e., the accessory genes).

7

Figure 2. Average sequence similarity (%) values for the various pairwise comparisons of
the five *Ceratocystis albifundus* genomes calculated using JSpecies (Richter and RossellóMóra 2009). The conserved proportions (%) of the genomes used for these comparisons are
indicated in parentheses. JSpecies artificially sectioned individual genomes into fragments
consisting of 100-1020 nucleotides, followed by pairwise comparisons using BLAST.
Average sequence similarity was estimated only for those fragments that aligned over > 70 %
of their entire lengths and were > 30 % similar.

15

16 Figure 3. Visualization of single nucleotide polymorphism (SNP) density, read depth, 17 distribution of genes associated with host interactions and movement of TEs, as well as TEs 18 and repeat regions along contig 7 of the Ceratocystis albifundus reference genome for isolate 19 CMW 4068. (A) The peaks at the top represent SNP distribution that was defined by non-20 overlapping 5000 bp intervals across the sequence and measuring the SNP content as the number of SNPs and read depth per interval using the R/Bioconductor package KaryoplotEr 21 22 (Gel and Serra 2017). The peaks at the bottom represent read depth calculated using a sliding window of 5000 bp. The positions of transposable elements and repeat regions are indicated 23 24 with dark brown vertical lines, putative CAZymes identified using dbCAN database (Yin et al. 2012) are indicated with yellow lines, putative peptidases identified using the MEROPS 25 26 (Rawlings et al. 2016) are indicated with purple lines, putative effectors identified using 27 EffectorP v1.0 (Sperschneider et al. 2016) are indicated with blue horizontal lines and genes 28 with significant similarity to those previously shown to be involved in the movement and activity of TEs that are present in the Gypsy DataBase (Llorens et al. 2011) are indicated with 29 30 green horizontal lines. Synteny breaks between the reference and query genomes were determined using SynChro (Drillon et al. 2014). (B) MAUVE visualisation of synteny 31 32 between the five Ceratocystis albifundus genomes. Pairwise alignments of genomes were generated using the MAUVE plugin implemented in Geneious v. 7 (Kearse et al. 2012). 33 34 Locally Collinear Blocks are marked with the same colour and connected by straight lines.

Figure 4. (A) Genome organization of the the five Ceratocystis albifundus genomes. 2 3 Conserved synteny blocks were defined between pairwise combinations of the five genomes 4 using Synchro (options: 0 3; 0 for all pairwise, and 2 for delta of RBH genes) (Drillon et al. 5 2014). This program defined orthology relationships between genes from different isolates on 6 the basis of bidirectional hits in a BLASTp comparison (reciprocal best hits). Different 7 colours are used to differentiate gene contents in different ancestral C. albifundus contigs. 8 The colour white indicated the absence of orthologs in the other isolate. Classes (B) and orders (C) of TEs identified in the five Ceratocystis albifundus genomes examined based on 9 the classification scheme of Wicker et al. (2007). The various elements were denoted as 10 follows: TIR = Terminal Inverted Repeats; MITE = Miniature Inverted Repeat; LINE = Long 11 12 Interspersed Nuclear Elements; LTR = Long Terminal Repeat; LARD = Large retrotransposon derivatives; TRIM = Terminal repeat transposons in miniature; SINE = Short 13 Interspersed Elements; and Unclassified = non-autonomous retrotransposons. Those in Class 14 I included LTR, LINES, SINES, LARD and TRIM), Class II included TIR and MITE, while 15 16 the unclassified TEs formed part of NoClass.

17

1

- 18
- 19
- 20

21

1 Supplementary Figures

2

3 Figure S1. Heatmap showing copy number and relationships among the genes included in 4 the non-shared set that had similarity to entries in the Kyoto Encyclopaedia of Genes and 5 Genomes (KEGG) ORTHOLOGY database (Table S5). K numbers, KEGG definitions and 6 pathways associated with each entry are indicated on the right. Gene copy numbers are 7 indicated according to the scale bar. For the heatmap, the relationships among the genomes 8 were inferred from the gene copy number data using Euclidean Correlation distances and 9 average linkage between the taxa using ClustVis (Metsalu and Vilo 2015). The data were 10 normalized using row centering and unit variance scaling (i.e., divides the values by standard 11 deviation so that each row has variance equal to one).

12

Figure S2. Putative *Ceratocystis albifundus* phospholipases identified by BLASTp searches
(E-value cut-off of ≤ 0.00001) against previously characterized phospholipases (PLA2,
KEY75421; PLB2, KEY77760; PLC2, XP_011319931). These sequences were aligned using
MAFFT v. 7 (http://mafft.cbrc.jp/alignment/server/), followed by phylogenetic analysis with
a neighbor-joining method in the software MEGA v. 7 (http://www.megasoftware.net).

18

Figure S3. Putative *Ceratocystis albifundus* catalases and catalase-peroxidases identified by
BLASTp searches (E-value cut-off of ≤ 0.00001) against previously characterized catalases
(CATA, accession number MG100061; CATB, MG06442) and catalase-peroxidases
(KATG1, A4R5S9; KATG2, A4QUT2). These sequences were aligned using MAFFT v. 7
(http://mafft.cbrc.jp/alignment/server/), followed by phylogenetic analysis with a neighborjoining method in the software MEGA v. 7 (http://www.megasoftware.net).

25

26 Figure S4. Putative Ceratocystis albifundus superoxide dismutases identified by BLASTp 27 searches (E-value cut-off of ≤ 0.00001) against previously characterized superoxide dismutases (SOD1, EFZ03762; SOD2, EFY99820; SOD3, EFY99375; SOD4, EFZ00595; 28 EFZ00365). 29 SOD5, These sequences aligned using MAFFT 7 were v. 30 (http://mafft.cbrc.jp/alignment/server/), followed by phylogenetic analysis with a neighborjoining method in the software MEGA v. 7 (http://www.megasoftware.net). 31

32

Figure S5. REVIGO (Supek et al. 2011) treemap summarizing GO biological process categories enriched in the accessory set. GO term enrichment in the accessory set

1 (Supplementary Table 9) was determined by performing a Fisher's exact test (two-sided; P <

- 2 0.05) using Blast2GO (Conesa et al. 2005).
- 3

4 Figure S6. Visualization of single nucleotide polymorphism (SNP) density, read depth, 5 distribution of genes associated with host interactions and movement of TEs, as well as TEs 6 and repeat regions along the 10 largest contigs of the Ceratocystis albifundus reference 7 genome for isolate CMW 4068. (A) The peaks at the top represent SNP distribution that was 8 defined by non-overlapping 5000 kb intervals across the sequence and measuring the SNP content as the number of SNPs and read depth per interval using the R/Bioconductor package 9 KaryoplotEr (Gel and Serra 2017). The peaks at the bottom represent read depth calculated 10 using a sliding window of 5000 kb. Transposable elements and repeat regions are indicated 11 with dark brown vertical lines, putative CAZymes identified using dbCAN database (Yin et 12 al. 2012) are indicated with yellow lines, putative peptidases identified using the MEROPS 13 (Rawlings et al. 2016) are indicated with purple lines, putative effectors identified using 14 15 EffectorP v1.0 (Sperschneider et al. 2016) are indicated with blue horizontal lines and genes with significant similarity to those previously shown to be involved in the movement and 16 activity of TEs that are present in the Gypsy DataBase (Llorens et al. 2011) are indicated with 17 green horizontal lines. (B) MAUVE visualisation of synteny between the 5 Ceratocystis 18 albifundus genomes. Pairwise alignments of genomes were generated using MAUVE pluging 19 plugin implemented in Geneious v. 7 (Kearse et al. 2012). Locally Collinear Blocks (LCBs) 20 are marked with the same colour and connected by straight lines. 21

22

23

24

1	Supplementary Tables
2	Supplementary Table 1. Sequence statistics for the Ceratocystis albifundus genomes.
3	
4	Supplementary Table 2. BUSCO results for the Ceratocystis albifundus genomes.
5	
6	Supplementary Table 3. Inventory of the genes present in the Ceratocystis albifundus
7	reference genome (Isolate CMW 4068).
8	
9	Supplementary Table 4. Inventory of the genes present in all of the five Ceratocystis
10	albifundus genomes (referred to as the core genes set), identified using reciprocal Basic Local
11	Alignment Search Tool (BLAST) searches.
12	
13	Supplementary Table 5. Inventory of the genes predicted to be present in two to four of the
14	five Ceratocystis albifundus genomes, identified using reciprocal BLASTp searches.
15	
16	Supplementary Table 6. Inventory of the genes unique to a specific Ceratocystis albifundus
17	assembly, identified using reciprocal BLASTp searches.
18	
19	Supplementary Table 7. A list of the genes present in all of the five Ceratocystis albifundus
20	genomes that were associated with pathways reconstructed using the Kyoto Encyclopaedia of
21	Genes and Genomes (KEGG) databases (http://www.genome.jp/kegg/).
22	
23	Supplementary Table 8. A list of the Ceratocystis albifundus accessory genes that had
24	significant similarity to entries in the Kyoto Encyclopaedia of Genes and Genomes (KEGG)
25	databases (http://www.genome.jp/kegg/).
26	
27	Supplementary Table 9. GO terms significantly enriched in the <i>Ceratocystis</i> accessory set.
28	Y'
29	Supplementary Table 10. The genes listed in the Pathogen-Host Interaction (PHI) database
30	(Winnenburg et al. 2008) predicted to be present in the Ceratocystis albifundus assemblies.
31	
32	Supplementary Table 11. Putative Ceratocystis albifundus peptidases identified using the
33	MEROPS database (http://merops.sanger.ac.uk; Rawlings et al. 2016).
34	

1	Supplementary Table 12. Putative Ceratocystis albifundus Carbohydrate-active enzymes
2	(CAZymes) identified and annotated using the online resource dbCAN (DataBase for
3	automated Carbohydrate-active enzyme ANnotation) (Yin et al. 2012).
4	
5	Supplementary Table 13. Transposable elements identified in the five Ceratocystis
6	albifundus accessory gene set identified using BLASTp against the Gypsy Database (GyDB)
7	of mobile genetic elements: release 2.0 (Llorens et al. 2011).
8	
9	Supplementary Table 14. Mapping of Illumina assemblies to the reference genome (Isolate
10	CMW 4068) using the LASTZ (Large-Scale Genome Alignment Tool) plugin implemented
11	in Geneious v. 7 (Kearse et al. 2012).
12	
13	Supplementary Table 15. Mapping results and single nucleotide polymorphism (SNP)
14	identified for each of isolates CMW 17620, CMW 17274, CMW 13980 and CMW 17274
15	using the reference genome (Isolate CMW 4068) and CLC Genomics Workbench v. 6.0.1.
16	
17	Supplementary Table 16. Likelihood scores and parameter estimates for the site-specific
18	selection models (Yang et al. 2000) evaluated in this study.
19	
20	Supplementary Table 17. Transposable elements identified in the Ceratocystis albifundus
21	reference genome (Isolate CMW 4068) identified and annotated using TEdenovo and the
22	TEannot pipelines from the REPET package (Flutre et al. 2011).
23	
24	

1 References

- 2 Altschul S.F., Madden T.L., Schäffer A.A., Zhang J., Zhang Z., Miller W., Lipman D.J.,
- 3 1997. Gapped BLAST and PSI-BLAST: a new generation of protein database search
 4 programs. *Nucleic Acids Res.* 25, 3389-3402.
- 5 Amselem J., Cuomo C.A., Van Kan J.A., Viaud M., Benito E.P., Couloux A., Coutinho P.M.,

De Vries R.P., Dyer P.S., Fillinger S., Fournier E., 2011. Genomic analysis of the
necrotrophic fungal pathogens *Sclerotinia sclerotiorum* and *Botrytis cinerea*. *PLOS Genet*.

- 8 7, e1002230.
- 9 Amselem J., Lebrun M.-H., Quesneville H., 2015. Whole genome comparative analysis of
 10 transposable elements provides new insight into mechanisms of their inactivation in fungal
 11 genomes. *BMC Genomics*. 16, 141.

12 Aylward J., Wingfield B.D., Dreyer L.L., Roets F., Wingfield M.J., Steenkamp E.T., 2017.

13 Contrasting carbon metabolism in saprotrophic and pathogenic microascalean fungi from

- 14 Protea trees. *Fungal Ecol.* 30, 88-100.
- Bao Z, Eddy S.R., 2002. Automated *de novo* identification of repeat sequence families in
 sequenced genomes. *Genome Res.* 12, 1269-1276.

Barnes I., Gaur A., Burgess T., Roux J., Wingfield B.D., Wingfield M.J., 2001. Microsatellite
 markers reflect intra-specific relationships between isolates of the vascular wilt pathogen
 Ceratocystis fimbriata. Mol. Plant Pathol. 2, 319-325.

20 Barnes I., Nakabonge G., Roux J., Wingfield B.D., Wingfield M.J., 2005. Comparison of

- populations of the wilt pathogen *Ceratocystis albifundus* in South Africa and Uganda. *Plant Pathol.* 54, 189-195.
- 23 Berlin K., Koren S., Chin C.S., Drake J.P., Landolin J.M., Phillippy A.M., 2015. Assembling

large genomes with single-molecule sequencing and locality-sensitive hashing. *Nat. Biotechnol.* 33, 623-630.

- Biely P., 2012. Microbial carbohydrate esterases deacetylating plant polysaccharides.
 Biotechnol. Adv. 30, 1575-1588.
- Biémont C., 2010. A brief history of the status of transposable elements: from junk DNA to
 major players in evolution. *Genetics*. 186, 1085-1093.
- 30 Bland N.D., Pinney J.W., Thomas J.E., Turner A.J., Isaac R.E., 2008. Bioinformatic analysis
- of the neprilysin (M13) family of peptidases reveals complex evolutionary and functional
 relationships. *BMC Evol. Biol.* 23, 1.
- Boetzer M, Pirovano W., 2014. SSPACE-LongRead: scaffolding bacterial draft genomes
 using long read sequence information. *BMC Bioinformatics*. 15, 211.

- Boetzer M., Henkel C.V., Jansen H.J., Butler D., Pirovano W., 2011. Scaffolding preassembled contigs using SSPACE. *Bioinformatics*. 27, 578-579.
- Boetzer M., Pirovano W., 2012. Toward almost closed genomes with GapFiller. *Genome Biol.* 13, R56.
- Böhne A., Brunet F., Galiana-Arnoux D., Schultheis C., Volff J.-N., 2008. Transposable
 elements as drivers of genomic and biological diversity in vertebrates. *Chromosome Res.*16, 203-215.
- 8 Braunsdorf C., Mailänder-Sánchez D., Schaller M., 2016. Fungal sensing of host
 9 environment. *Cell. Microbiol.* 18,1188-2000.
- Cantarel B.L., Coutinho P.M., Rancurel C., Bernard T., Lombard V., Henrissat B., 2009. The
 Carbohydrate-Active EnZymes database (CAZy): an expert resource for glycogenomics.

12 *Nucleic Acids Res.* 37, D233-238.

- 13 Casacuberta J.M., Santiago N., 2003. Plant LTR-retrotransposons and MITEs: control of
- transposition and impact on the evolution of plant genes and genomes. *Gene*. 311, 1-11.
- Chiapello H., Mallet L., Guérin C., Aguileta G., Amselem J., Kroj T., Ortega-Abboud E.,
 Lebrun M.H., Henrissat B., Gendrault A., Rodolphe F., 2015. Deciphering genome
 content and evolutionary relationships of isolates from the fungus *Magnaporthe oryzae*attacking different host plants. *Genome Biol. Evol.* 7, 2896-2912.
- 19 Conesa A., Götz S., García-Gómez J.M., Terol J., Talón M., Robles M., 2005. Blast2GO: a
- universal tool for annotation, visualization and analysis in functional genomics research. *Bioinformatics*. 18, 3674-3676.
- Croll D., McDonald B.A., 2012. The accessory genome as a cradle for adaptive evolution in
 pathogens. *PLOS Pathog.* 8, e1002608.
- 24 Daboussi M.-J., 1997. Fungal transposable elements and genome evolution. *Genetica*. 100,
 25 253-260.
- Daboussi M.-J., Capy P., 2003. Transposable elements in filamentous fungi. *Annu. Rev. Microbiol.* 57, 275–299.
- Daboussi M.J., 1996. Fungal transposable elements: generators of diversity and genetic tools. *J. Genet.* 75, 325-339.
- Darling A.C., Mau B., Blattner F.R., Perna N.T., 2004. Mauve: multiple alignment of
 conserved genomic sequence with rearrangements. *Genome Res.* 14, 1394-1403.
- 32 Desjardins A.E., Hohn T.M., 1997. Mycotoxins in plant pathogenesis. Mol. Plant Microbe
- 33 *Interact.* 10, 147-152.

- 1 Dong S., Raffaele S., Kamoun S., 2015. The two-speed genomes of filamentous pathogens:
- 2 waltz with plants. *Curr. Opin. Genet. Dev.* 35, 57-65.
- 3 Dopman E., Hartl D., 2007. A portrait of copy-number polymorphism in *Drosophila*4 *melanogaster. Proc. Natl. Acad. Sci. USA.* 104, 19920-19925.
- 5 Drillon G., Carbone A., Fischer G., 2014. SynChro: a fast and easy tool to reconstruct and
 6 visualize synteny blocks along eukaryotic chromosomes. *PLOS ONE*. 9, e92621.
- 7 Edgar R.C., Myers E.W., 2003. PILER: Identification and classification of genomic repeats.

8 *Bioinformatics*. 21, 152-158.

- 9 Ellinghaus D., Kurtz S., Willhoeft U., 2008. LTRharvest, an efficient and flexible software
 10 for de novo detection of LTR retrotransposons. *BMC Bioinformatics*. 9, 18.
- Faino L., Seidl M.F., Shi-Kunne X., Pauper M., van den Berg G.C., Wittenberg A.H.,
 Thomma B.P., 2016. Transposons passively and actively contribute to evolution of the
 two-speed genome of a fungal pathogen. *Genome Res.* 26, 1091-1100.
- Finn, R.D., Clements, J., Eddy, S.R., 2011. HMMER web server: interactive sequence
 similarity searching. *Nucleic Acids Res.* 39, W29–W37.
- Flutre T., Duprat E., Feuillet C., Quesneville H., 2011. Considering transposable element
 diversification in *de novo* annotation approaches. *PLOS ONE* 6, e16526.
- 18 Fortna A., Kim Y., MacLaren E., Marshall K., Hahn G., Meltesen L., Brenton M., Hink R.,
- Burgers S., Hernandez-Boussard T., Karimpour-Fard A., 2004. Lineage-specific gene
 duplication and loss in human and great ape evolution. *PLOS Biol.* 2, e207.
- Fouché S., Plissonneau C., Croll D., 2018. The birth and death of effectors in rapidly
 evolving filamentous pathogen genomes. *Curr. Opin. Microbiol.* 46, 34-42.
- 23 Fudal I., Ross S., Brun H., Besnard A.-L., Ermel M., Kuhn M.-L., Balesdent M.-H., Rouxel
- 24 T., 2009. Repeat-induced point mutation (RIP) as an alternative mechanism of evolution
- toward virulence in *Leptosphaeria maculans*. *Mol. Plant Microbe Interact*. 22, 932-941.
- 26 Fulnečková J., Ševčíková T., Fajkus J., Lukešova A., Lukeš M., Vlček Č., Lang B.F., Kim E.,
- Eliaš M., Sýkorova E., 2013. A broad phylogenetic survey unveils the diversity and
 evolution of telomeres in eukaryotes. *Genome Biol. Evol.* 5, 468-483.
- Gel B., Serra E., 2017. karyoploteR: An R/Bioconductor package to plot customizable linear
 genomes displaying arbitrary data. *BioRxiv*. 122838.
- 31 Ghannoum M., 2000. Potential role of phospholipases in virulence and fungal pathogenesis.
- **32** *Clinical Microbiol. Rev.* 13, 122-143.
- 33 Gish W., States D.J., 1993. Identification of protein coding regions by database similarity
- 34 search. *Nat. Genet.* 3, 266-272.

1	Gladieux P., Ropars J., Badouin H., Branca A., Aguileta G., Vienne D.M., Rodríguez de la
2	Vega R.C., Branco S., Giraud T., 2014. Fungal evolutionary genomics provides insight
3	into the mechanisms of adaptive divergence in eukaryotes. Mol. Ecol. 23, 753-773.
4	Goodwin S.B., M'Barek S.B., Dhillon B., Wittenberg A.H., Crane C.F., Hane J.K., Foster
5	A.J., Van der Lee T.A., Grimwood J., Aerts A., Antoniw J., 2011. Finished genome of the
6	fungal wheat pathogen Mycosphaerella graminicola reveals dispensome structure,
7	chromosome plasticity, and stealth pathogenesis. PLOS Genet. 7, e1002070.
8	Goris J., Konstantinidis K.T., Klappenbach J.A., Coenye T., Vandamme P., Tiedje J.M.,
9	2007. DNA-DNA hybridization values and their relationship to whole-genome sequence
10	similarities. Int. J. Syst. Evol. Microbiol. 57, 81-91.
11	Grandaubert J., Lowe R.G., Soyer J.L., Schoch C.L., Van de Wouw A.P., Fudal I., Robbertse
12	B., Lapalu N., Links M.G., Ollivier B., Linglin J., 2014. Transposable element-assisted
13	evolution and adaptation to host plant within the Leptosphaeria maculans-Leptosphaeria
14	biglobosa species complex of fungal pathogens. BMC Genomics. 15, 1.
15	Hall T., 2011. BioEdit: An important software for molecular biology. GERF Bull. Bioscience.
16	2, 60-61.
17	Hardison R.C., 2003. Comparative genomics. PLOS Biol. 1, 156-160.
18	Harris R.S., 2007. Improved pairwise alignment of genomic DNA, The Pennsylvania State
19	University.
20	Hartl L., Zach S., Seidl-Seiboth V., 2012. Fungal chitinases: diversity, mechanistic properties
21	and biotechnological potential. Appl. Microbiol. Biotechnol. 93, 533-543.
22	Hatta R., Ito K., Hosaki Y., Tanaka T., Tanaka A., Yamamoto M., Akimitsu K., Tsuge T.,
23	2002. A conditionally dispensable chromosome controls host-specific pathogenicity in the
24	fungal plant pathogen Alternaria alternata. Genetics. 161, 59-70.
25	Heath R.N., Wingfield M.J., Wingfield B.D., Meke G., Mbaga A., Roux J., 2009.
26	Ceratocystis species on Acacia mearnsii and Eucalyptus spp. in eastern and southern
27	Africa including six new species. Fungal Div. 34, 41-67.
28	Kanehisa M., Sato Y., Morishima K., 2016. BlastKOALA and GhostKOALA: KEGG tools
29	for functional characterization of genome and metagenome sequences. J. Mol. Biol. 428,
30	726-731.
31	Kearse M., Moir R., Wilson A., Stones-Havas S., Cheung M., Sturrock S., Buxton S., Cooper
32	A., Markowitz S., Duran C., Thierer T., 2012. Geneious basic: an integrated and
33	extendable desktop software platform for the organization and analysis of sequence data.
34	Bioinformatics. 28, 1647-1649.

1	Kidwell M.G., Lisch D.R., 2000. Transposable elements and host genome evolution. Trends
2	Ecol. Evol. 15, 95-99.
3	Labbe J., Murat C., Morin E., Tuskan G.A., Le Tacon F., Martin F., 2012. Characterization of
4	transposable elements in the ectomycorrhizal fungus Laccaria bicolor. PLOS ONE, 7,
5	e40197.
6	Laurie J.D., Ali S., Linning R., Mannhaupt G., Wong P., Güldener U., Münsterkötter M.,
7	Moore R., Kahmann R., Bakkeren G., Schirawski J., 2012. Genome comparison of barley
8	and maize smut fungi reveals targeted loss of RNA silencing components and species-
9	specific presence of transposable elements. <i>Plant Cell</i> . 24, 1733-1745.
10	Leclair S., Ansan-Melayah D., Rouxel T., Balesdent M., 1996. Meiotic behavior of the
11	minichromosome in the phytopathogenic ascomycete Leptosphaeria maculans. Curr.
12	Genetics. 30, 541-548.
13	Lee D.H., Roux J., Wingfield B.D., Barnes I., Mostert L., Wingfield M.J., 2016. The genetic
14	landscape of Ceratocystis albifundus populations in South Africa reveals a recent fungal
15	introduction event. Fungal Biol. 120, 690-700.
16	Lee M.J., Sheppard D.C., 2016. Recent advances in the understanding of the Aspergillus
17	fumigatus cell wall. J. Microbiol. 54, 232-242.
18	Levasseur A., Drula E., Lombard V., Coutinho P.M., Henrissat B., 2013. Expansion of the
19	enzymatic repertoire of the CAZy database to integrate auxiliary redox enzymes.
20	Biotechnol. Biofuels. 6, 1.
21	Li F., Shi HQ., Ying SH., Feng MG., 2015. Distinct contributions of one Fe- and two
22	Cu/Zn-cofactored superoxide dismutases to antioxidation, UV tolerance and virulence of
23	Beauveria bassiana. Fungal Genet. Biol. 81, 160–171.
24	Llorens C., Futami R., Covelli L., Dominguez-Escriba L., Viu J.M., Tamarit D., Aguilar-
25	Rodriguez J., Vicente-Ripolles M., Fuster G., Bernet G.P., Maumus F., 2011. The Gypsy
26	Database (GyDB) of Mobile Genetic Elements: Release 2.0. Nucleic Acids Res. 39 (suppl
27	1), D70-D74.
28	Ma L.J., Van Der Does H.C., Borkovich K.A., Coleman J.J., Daboussi M.J., Di Pietro A.,
29	Dufresne M., Freitag M., Grabherr M., Henrissat B., Houterman P.M., 2010. Comparative
30	genomics reveals mobile pathogenicity chromosomes in Fusarium. Nature. 464, 367-373.
31	Ma L.J., Geiser D.M., Proctor R.H., Rooney A.P., O Donnell K., Trail F., Gardiner D.M.,
32	Manners J.M., Kazan K., 2013. Fusarium pathogenomics. Annu. Rev. Microbiol. 67, 399-
33	416.

1	Manning V.A., Pandelova I., Dhillon B., Wilhelm L.J., Goodwin S.B., Berlin A.M., Figueroa
2	M., Freitag M., Hane J.K., Henrissat B., Holman W.H., 2013. Comparative genomics of a
3	plant-pathogenic fungus, Pyrenophora tritici-repentis, reveals transduplication and the
4	impact of repeat elements on pathogenicity and population divergence. G3 Genes Genom.
5	Genet. 3, 41-63.
6	Metsalu T., Vilo J., 2015. Clustvis: a web tool for visualizing clustering of multivariate data
7	using Principal Component Analysis and heatmap. Nucleic Acids Res. 43, W566-W70.
8	Mukherjee P.K., Horwitz B.A., Kenerley C.M., 2012. Secondary metabolism in Trichoderma
9	- a genomic perspective. <i>Microbiol</i> . 158, 35-45.
10	Novikova O., 2009. Chromodomains and LTR retrotransposons in plants. Commun. Integr.
11	Biol. 2, 158-162.
12	Ohm R.A., Feau N., Henrissat B., Schoch C.L., Horwitz B.A., Barry K.W., Condon B.J.,
13	Copeland A.C., Dhillon B., Glaser F., Hesse C.N., 2012. Diverse lifestyles and strategies
14	of plant pathogenesis encoded in the genomes of eighteen Dothideomycetes fungi. PLOS
15	Pathog. 8, e1003037.
16	Paris S., Wysong D., Debeaupuis JP., Shibuya K., Philippe B., Diamond, R.D., Latge, JP.,
17	2003. Catalases of Aspergillus fumigatus. Infect. Immun. 71, 3551-3562.
18	Petersen T.N., Brunak S., von Heijne G., Nielsen H., 2011. SignalP 4.0: discriminating signal
19	peptides from transmembrane regions. Nat. Methods. 8, 785-786.
20	Plissonneau C., Benevenuto J., Mohd-Assaad N., Fouché S., Hartmann F.E., Croll D., 2017.
21	Using population and comparative genomics to understand the genetic basis of effector-
22	driven fungal pathogen evolution. Front. Plant Sci. 8, 119.
23	Plissonneau C., Hartmann F.E., Croll D., 2018. Pangenome analyses of the wheat pathogen
24	Zymoseptoria tritici reveal the structural basis of a highly plastic eukaryotic genome. BMC
25	<i>Biol.</i> 16, 5.
26	Plissonneau C., Stürchler A., Croll D., 2016. The evolution of orphan regions in genomes of
27	a fungal pathogen of wheat. MBio. 7, e01231-16.
28	Quesneville H., Nouaud D., Anxolabéhère D., 2003. Detection of new transposable element
29	families in Drosophila melanogaster and Anopheles gambiae genomes. J. Mol. Evol. S50-
30	S59.
31	Raffaele S., Kamoun S., 2012. Genome evolution in filamentous plant pathogens: why bigger
32	can be better. Nat. Rev. Microbiol. 10, 417-430.
33	Rawlings N.D., Barrett A.J., Finn R.D., 2016. Twenty years of the MEROPS database of
34	proteolytic enzymes, their substrates and inhibitors. Nucleic Acids Res. 44, D343-D350.

- Richter M., Rosselló-Móra R., 2009. Shifting the genomic gold standard for the prokaryotic
 species definition. *Proc. Natl. Acad. Sci. USA*. 106, 19126-19131.
- Roux J., Harrington T.C., Steimel J.P., Wingfield M.J., 2001. Genetic variation in the wattle
 wilt pathogen *Ceratocystis albifundus*. *Mycoscience*. 42, 327-332.
- Roux J., Heath R.N., Labuschagne L., Nkuekam G.K., Wingfield M.J., 2007. Occurrence of
 the wattle wilt pathogen, *Ceratocystis albifundus* on native South African trees. *Forerst Pathol.* 37, 292-302.
- Roux J., Dunlop R., Wingfield M.J., 1999. Susceptibility of elite *Acacia mearnsii* families to
 Ceratocystis wilt in South Africa. *J. Forest Res.* 4, 187-190.
- 10 Roux J., Meke G., Kanyi B., Mwangi L., Mbaga A., Hunter G.C., Nakabonge G., Heath R.N.,
- 11 Wingfield M.J., 2005. Diseases of plantation forestry trees in eastern and Southern Africa.
- 12 S. Afr. J. Sci. 101, 409-413.
- 13 Shi X., Faino L., van den Berg G., Thomma B., Seidl M., 2018. Evolution within the fungal
- genus *Verticillium* is characterized by chromosomal rearrangement and gene loss. *Environ. Microbiol.* 20, 1362-73.
- Simão F.A., Waterhouse R.M., Ioannidis P., Kriventseva E.V., Zdobnov E.M., 2015.
 BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*. 31, 3210-3212.
- 19 Sperschneider J., Gardiner D.M., Dodds P.N., Tini F., Covarelli L., Singh K.B., Manners
- J.M., Taylor J.M., 2016. EffectorP: predicting fungal effector proteins from secretomes
 using machine learning. *New Phytol.* 210, 743-761.
- Stanke M., Steinkamp R., Waack S., Morgenstern B., 2004. AUGUSTUS: A web server for
 gene finding in Eukaryotes. *Nucleic Acids Res.* 32, W309-W312.
- Steenkamp E.T., Wingfield M.J., McTaggart A.R., Wingfield B.D., 2018. Fungal species and
 their boundaries matter Definitions, mechanisms and practical implications. *Fungal Biol.*
- **26** *Rev.* 32, 104-116.
- Stergiopoulos I., de Wit P.J., 2009. Fungal effector proteins. *Annu. Rev. Phytopathol.* 8, 233263.
- Supek F., Bosnjak M., Skunca N., Smuc T., 2011. REVIGO summarizes and visualizes long
 lists of gene ontology terms. *PLOS ONE*. 6, e21800.
- 31 Tanabe S., Ishii-Minami N., Saitoh K.-I., Otake Y., Kaku H., Shibuya N., NishizawaY.,
- 32 Minami E., 2011. The role of catalase-peroxidase secreted by *Magnaporthe oryzae* during 32 $M_{1} = \frac{1}{2} \int \frac$
- arly infection of rice cells. *Mol. Plant Microbe Interact.* 24, 163-171.

1	Tzeng T.H., Lyngholm L.K., Ford C.F., Bronson C.R., 1992. A restriction fragment length
2	polymorphism map and electrophoretic karyotype of the fungal maize pathogen
3	Cochliobolus heterostrophus. Genetics. 130, 81-96.
4	van der Nest M.A., Beirn L.A., Crouch J.A., Demers J.E., De Beer Z.W., De Vos L., Gordon
5	T.R., Moncalvo JM., Naidoo K., Sanchez-Ramirez S., Roodt D., 2014a. Draft genomes
6	of Amanita jacksonii, Ceratocystis albifundus, Fusarium circinatum, Huntiella omanensis,
7	Leptographium procerum, Rutstroemia sydowiana, and Sclerotinia echinophila. IMA
8	Fungus. 5, 473-486.
9	van der Nest M.A., Bihon W., De Vos L., Naidoo K., Roodt D., Rubagotti E., Slippers B.,
10	Steenkamp E.T., Wilken P.M., Wilson A., Wingfield M.J., Wingfield B.D., 2014b. Draft
11	genome sequences of Diplodia sapinea, Ceratocystis manginecans, and Ceratocystis
12	moniliformis. IMA Fungus. 5, 135-140.
13	van der Nest M.A., Steenkamp E.T., McTaggart A.R., Trollip C., Godlonton T., Sauerman
14	E., Roodt D., Naidoo K., Coetzee M.P.A., Wilken P.M., 2015. Saprophytic and pathogenic
15	fungi in the Ceratocystidaceae differ in their ability to metabolize plant-derived sucrose.
16	BMC Evol. Biol. 15, 1.
17	van Wyk S., Wingfield B.D., De Vos L., Santana Q., Van der Merwe N., Steenkamp E.T.,
18	2018. Multiple independent origins for a subtelomeric locus associated with growth rate in
19	Fusarium circinatum. IMA Fungus. 9, 27–36.
20	Vanheule A., Audenaert K., Warris S., van de Geest H., Schijlen E., Höfte M., De Saeger S.,
21	Haesaert G., Waalwijk C., van der Lee T., 2016. Living apart together: crosstalk between
22	the core and supernumerary genomes in a fungal plant pathogen. BMC Genomics. 17, 670.
23	Vanrobays E., Thomas M., Tatout C., 2018. Heterochromatin positioning and nuclear
24	architecture. Annu. Rev. Plant Biol. 46, 157-190.
25	Walker B.J., Abeel T., Shea T., Priest M., Abouelliel A., Sakthikumar S., Cuomo C., Zeng
26	Q., Wortman J., Young S.K., Earl A.M., 2014. Pilon: An integrated tool for
27	comprehensive microbial variant detection and genome assembly improvement. PLOS
28	ONE. 9, e112963.
29	Wicker T., Sabot F., Hua-Van A., Bennetzen J.L., Capy P., Chalhoub B., Flavell A., Leroy
30	P., Morgante M., Panaud O., 2007. A unified classification system for eukaryotic
31	transposable elements. Nat. Rev. Genet. 8, 973-982.
32	Wilhelm M., Wilhelm F.X., 2001. Reverse transcription of retroviruses and LTR
33	retrotransposons. CMLS-Cell. Mol. Life S. 58, 1246-1262.

- Wilken P.M., Steenkamp E.T., Wingfield M.J., De Beer Z.W., Wingfield B.D., 2013.
 Ceratocystis fimbriata: Draft nuclear genome sequence for the plant pathogen,
- 3 *Ceratocystis fimbriata. IMA Fungus.* 4, 357-358.

Wingfield B.D., Barnes I., de Beer Z.W., De Vos L., Duong T.A., Kanzi A.M., Naidoo K.,
Nguyen H.D., Santana Q.C., Sayari M., Seifert K.A., 2015. IMA Genome-F 5: Draft
genome sequences of *Ceratocystis eucalypticola, Chrysoporthe cubensis, C. deuterocubensis, Davidsoniella virescens, Fusarium temperatum, Graphilbum fragrans, Penicillium nordicum,* and *Thielaviopsis musarum. IMA Fungus.* 6, 493-506.

- 9 Wingfield B.D., Ambler J.M., Coetzee M., De Beer Z.W., Duong T.A., Joubert F.,
 10 Hammerbacher A., McTaggart A.R., Naidoo K., Nguyen H.D., Ponomareva E., 2016a.
 11 Draft genome sequences of Armillaria fuscipes, Ceratocystiopsis minuta, Ceratocystis
 12 adiposa, Endoconidiophora laricicola, E. polonica and Penicillium freii. IMA Fungus 7,
 13 217-227.
- Wingfield B.D., Duong T.A., Hammerbacher A., van der Nest M.A., Wilson A., Chang R.,
 De Beer W., Steenkamp E.T., Wilken M.P., Naidoo K., Wingfield M.J., 2016b. Draft
 genome sequences for *Ceratocystis fagacearum*, *C. harringtonii*, *Grosmannia penicillata*,
 and *Huntiella bhutanensis*. *IMA Fungus* 7, 317-323.
- Winnenburg R., Urban M., Beacham A., Baldwin T.K., Holland S., Lindeberg M., Hansen
 H., Rawlings C., Hammond-Kosack K.E., Köhler J., 2008. PHI-base update: additions to
- 20 the pathogen-host interaction database. *Nucleic Acids Res.* 36 (suppl 1), D572-D576.
- Wittenberg A.H.J., Van der Lee T.A.J., Schouten H.J., 2009. Meiosis drives extraordinary
 genome plasticity in the haploid fungal plant pathogen *Mycosphaerella graminicola*. *PLOS ONE*. 4, 1-37.
- Yang Z., Nielson R., 2002. Codon-substitution models for detecting molecular adaptation at
 individual sites along specific lineages. *Mol. Biol. Evol.* 19, 908-917.
- Yang Z., Nielson R., Goldman N., Pedersen A., 2000. Codon-substitution models for
 heterogeneous selection pressure at amino acid sites. *Genetics*. 155, 431-449.
- Yin Y., Mao X., Yang J., Chen X., Mao F., Xu Y., 2012. dbCAN: a web resource for
 automated carbohydrate-active enzyme annotation. *Nucleic Acids Res.* 40, W445-W451.
- 30 Yue J.X., Li J., Aigrain L., Hallin J., Persson K., Oliver K., Bergström A., Coupland P.,
- Warringer J., Lagomarsino M.C., Fischer G., 2017. Contrasting evolutionary genome
 dynamics between domesticated and wild yeasts. *Nat. Genet.* 49, 913.
- Zerbino D.R., 2010. Using the Velvet *de novo* assembler for short-read sequencing
 technologies. *Curr. Protoc. Bioinformatics*. 11.5.1-11.5.12.

- 1 Zerbino D.R, Birney E., 2008. Velvet: Algorithms for de novo short read assembly using de
- 2 Bruijn graphs. *Genome Res.* 18, 821-829.
- 3 Zerillo M.M., Adhikari B.N., Hamilton J.P., Buell C.R., Levesque C.A., Tisserat N., 2013.
- Carbohydrate-active enzymes in *Pythium* and their role in plant cell wall and storage
 polysaccharide degradation. *PLOS ONE*. 8, e72572.
- 6 Zhao Z., Liu H., Wang C., Xu J.R., 2013. Comparative analysis of fungal genomes reveals
 7 different plant cell wall degrading capacity in fungi. *BMC Genomics*. 14, 1.
- 8 Zhou Y., Cahan SH., 2012. A novel family of Terminal-Repeat Retrotransposon in Miniature
- 9 (TRIM) in the genome of the red harvester Ant, *Pogonomyrmex barbatus*. *PLOS ONE*. 7,
 10 e53401.
- 11
- 12

Isolate number ^b	Geographical origin	Host	Collector(s)
CMW 4068	KwaZulu Natal, RSA	Acacia mearnsii (Fabaceae)	J. Roux
CMW 24685	Kenya	Acacia mearnsii (Fabaceae)	R.N. Heath & J. Roux
CMW 13980	Zambia	Parinari curatellifolia (Chrysobalanaceae)	J. Roux
CMW 17274	Gauteng, RSA	Faurea saligna (Proteaceae)	J. Roux
CMW 17620	Kruger National Park, RSA	Terminalia serecia (Combretaceae)	J. Roux

Table 1. Information about the *Ceratocystis albifundus* isolates used in this study^a.

^aThe genome assembly for isolate CMW17620 was available from a previous study (GenBank accession number: JSSU00000000; van der Nest et al. 2014), while those for the remaining isolates were determined here.

^bCMW: Culture collection of the Forestry and Agricultural Biotechnology Institute (FABI), University of Pretoria.

CER

Isolate number ^a	Size (Mp)	Nr. of large contigs	N50	N80	Largest contig	Nr. of genes ^d	Gene density (genes/Mb)
Illumina ^b							
CMW 4068	27.05	1 003	50 666	23 438	224 223	6 695	248
CMW 24685	26.97	1 072	48 267	21 815	289 214	6 710	249
CMW 13980	27.20	818	68 699	30 671	357 115	6 759	249
CMW 17274	26.68	2 122	23 089	9 324	153 182	6 699	251
CMW 17620	27.33	939	42 183	18 306	274 425	6 544	239
PacBio ^c							
CMW 4068	28.36	16	2 308 174	1 271 756	4 074 369	7 103	250

Table 2. Genome statistics for the five Ceratocystis albifundus isolates used in this study.

^aCMW: Culture collection of the Forestry and Agricultural Biotechnology Institute (FABI), University of Pretoria.

^bThe genome assembly for isolate CMW17620 was available from a previous study (GenBank accession number: JSSU00000000; van der Nest

et al. 2014), while those for the remaining isolates were determined here.

^cA high-quality reference genome was produced for isolate CMW 4068 using the PacBio and Illumina HiSeq sequencing platforms. The reference genome consisted of 16 contigs (> 200 000 bases).

^dORFs was predicted with the *de novo* gene prediction software AUGUSTUS, using the gene models of *Fusarium graminearum* (Stanke et al. 2004).

Table 3. KEGG functional classifications^a for the core genes shared by the five *Ceratocystis albifundus* isolates.

KEGG categories	Number of shared genes associated with pathway ^b			
Metabolism				
Carbohydrate metabolism	229			
Energy metabolism	122			
Lipid metabolism	128			
Nucleotide metabolism	130			
Amino acid metabolism	241			
Glycan biosynthesis and metabolism	76			
Metabolism of cofactors and vitamins	106			
Metabolism of terpenoids and polyketides	27			
Biosynthesis of other secondary metabolites	31			
Xenobiotics biodegradation and metabolism	46			
Genetic Information Processing				
Transcription	130			
Translation	298			
Folding, sorting and degradation	226			
Replication and repair	143			
Environmental Information Processing				
Membrane transport	5			
Signal transduction	385			
Cellular Processes				
Transport and catabolism	183			
Cell growth and death	230			
Cellular community	45			

^aFor these classifications, the predicted proteins were assigned to pathways using the Kyoto Encyclopaedia of Genes and Genomes (KEGG; http://www.genome.jp/kegg/) database and the GhostKoala mapping tool (Kanehisa et al. 2016).

^bThe analysis was done using the protein sequences for the 6241 genes shared by all five of the examined isolates.



		Target genome				
		CMW13980	CMW24685	CMW4068	CMW17274	CMW17620
	CMW13980	-	98,02 (95,73)	97,93 (95,44)	96,92 (94,98)	96,52 (85,06)
Query genome	CMW24685	97,57 (96,69)	-	98,63 (96,64)	97,90 (96,62)	97,15 (86,04)
	CMW4068	97,04 (96,16)	97,41 (96,48)	-	97,67 (95,59)	97,15 (86,10)
	CMW17274	97,26 (93,95)	98,40 (94,72)	98,00 (94,01)	-	96,89 (85,00)
	CMW17620	96,00 (85,04)	96,87 (85,30)	96,76 (85,29)	96,72 (85,91)	-





В









