

# First steps towards mitochondrial pan-genomics: Detailed analysis of *Fusarium graminearum* mitogenomes (#28541)

1

First submission

## Editor guidance

Please submit by **17 Aug 2018** for the benefit of the authors (and your \$200 publishing discount).



### Structure and Criteria

Please read the 'Structure and Criteria' page for general guidance.



### Custom checks

Make sure you include the custom checks shown below, in your review.



### Raw data check

Review the raw data. Download from the location [described by the author](#).



### Image check

Check that figures and images have not been inappropriately manipulated.

Privacy reminder: If uploading an annotated PDF, remove identifiable information to remain anonymous.

## Files

Download and review all files from the [materials page](#).

2 Figure file(s)

2 Latex file(s)

1 Table file(s)

1 Other file(s)

## ! Custom checks

### DNA data checks

! Have you checked the authors [data deposition statement](#)?

! Can you access the deposited data?

! Has the data been deposited correctly?

! Is the deposition information noted in the manuscript?



## Structure your review

The review form is divided into 5 sections. Please consider these when composing your review:

1. **BASIC REPORTING**
2. **EXPERIMENTAL DESIGN**
3. **VALIDITY OF THE FINDINGS**
4. General comments
5. Confidential notes to the editor






 You can also annotate this PDF and upload it as part of your review

When ready [submit online](#).





## Editorial Criteria

Use these criteria points to structure your review. The full detailed editorial criteria is on your [guidance page](#).





### BASIC REPORTING

-  Clear, unambiguous, professional English language used throughout.
-  Intro & background to show context. Literature well referenced & relevant.
-  Structure conforms to [PeerJ standards](#), discipline norm, or improved for clarity.
-  Figures are relevant, high quality, well labelled & described.
-  Raw data supplied (see [PeerJ policy](#)).

### EXPERIMENTAL DESIGN

-  Original primary research within [Scope of the journal](#).
-  Research question well defined, relevant & meaningful. It is stated how the research fills an identified knowledge gap.
-  Rigorous investigation performed to a high technical & ethical standard.
-  Methods described with sufficient detail & information to replicate.

### VALIDITY OF THE FINDINGS

-  Impact and novelty not assessed. Negative/inconclusive results accepted. *Meaningful* replication encouraged where rationale & benefit to literature is clearly stated.
-  Speculation is welcome, but should be identified as such.
-  Conclusions are well stated, linked to original research question & limited to supporting results.
-  Data is robust, statistically sound, & controlled.



The best reviewers use these techniques

## Tip

## Example

**Support criticisms with evidence from the text or from other sources**

*Smith et al (J of Methodology, 2005, V3, pp 123) have shown that the analysis you use in Lines 241-250 is not the most appropriate for this situation. Please explain why you used this method.*

**Give specific suggestions on how to improve the manuscript**

*Your introduction needs more detail. I suggest that you improve the description at lines 57- 86 to provide more justification for your study (specifically, you should expand upon the knowledge gap being filled).*

**Comment on language and grammar issues**

*The English language should be improved to ensure that an international audience can clearly understand your text. Some examples where the language could be improved include lines 23, 77, 121, 128 - the current phrasing makes comprehension difficult.*

**Organize by importance of the issues, and number your points**

1. Your most important issue
2. The next most important item
3. ...
4. The least important points

**Please provide constructive criticism, and avoid personal opinions**

*I thank you for providing the raw data, however your supplemental files need more descriptive metadata identifiers to be useful to future readers. Although your results are compelling, the data analysis should be improved in the following ways: AA, BB, CC*

**Comment on strengths (as well as weaknesses) of the manuscript**

*I commend the authors for their extensive data set, compiled over many years of detailed fieldwork. In addition, the manuscript is clearly written in professional, unambiguous language. If there is a weakness, it is in the statistical analysis (as I have noted above) which should be improved upon before Acceptance.*

# First steps towards mitochondrial pan-genomics: Detailed analysis of *Fusarium graminearum* mitogenomes

Balázs Brankovics <sup>Corresp., 1,2,3</sup>, Tomasz Kulik <sup>4</sup>, Jakub Sawicki <sup>4</sup>, Katarzyna Bilka <sup>4</sup>, Hao Zhang <sup>5</sup>, G Sybren de Hoog <sup>2,3</sup>, Theo AJ van der Lee <sup>1</sup>, Cees Waalwijk <sup>1</sup>, Anne D van Diepeningen <sup>1,2</sup>

<sup>1</sup> Wageningen Plant Research, Wageningen University & Research, Wageningen, Netherlands

<sup>2</sup> Westerdijk Fungal Biodiversity Institute, Utrecht, Netherlands

<sup>3</sup> Institute for Biodiversity and Ecosystem Dynamics, University of Amsterdam, Amsterdam, Netherlands

<sup>4</sup> Department of Botany and Nature Protection, University of Warmia and Mazury, Olsztyn, Poland

<sup>5</sup> State Key Laboratory for Biology of Plant Diseases and Insect Pests, Institute of Plant Protection, Chinese Academy of Agriculture Sciences, Beijing, China P.R.

Corresponding Author: Balázs Brankovics

Email address: balazs.brankovics@wur.nl

There is a gradual shift from representing a species' genome by a single reference genome sequence to a pan-genome representation. Pan-genomes are the abstract representations of the genomes of all the strains that are present in the population or species. In this study, we employed a pan-genomic approach to analyze the intraspecific mitochondrial genome diversity of *Fusarium graminearum*. We present an improved reference mitochondrial genome for *F. graminearum* with an intron-exon annotation that was verified using RNA-seq data. Each of the 24 studied isolates had a distinct mitochondrial sequence. Length variation in the *F. graminearum* mitogenome was found to be largely due to variation of intron regions (99.98%). The "intronless" mitogenome length was found to be quite stable and could be informative when comparing species. The coding regions showed high conservation, while the variability of intergenic regions was highest. However, the most important variable parts are the intron regions, because they contain approximately half of the variable sites, make up more than half of the mitogenome, and show presence/absence variation. Furthermore, our analyses show that the mitogenome of *F. graminearum* is recombining, as was previously shown in *F. oxysporum*, indicating that mitogenome recombination is a common phenomenon in *Fusarium*. The majority of mitochondrial introns in *F. graminearum* belongs to group I introns, which are associated with homing endonuclease genes (HEGs). Mitochondrial introns containing HE genes may spread within populations through homing, where the endonuclease recognizes and cleaves the recognition site in the target gene. After cleavage of the "host" gene, it is replaced by the gene copy containing the intron with HEG. We propose to use introns unique to a population for tracking the spread of the given population, because introns can spread through vertical inheritance, recombination as well

as via horizontal transfer. We demonstrated how pooled sequencing of strains can be used for mining mitogenome data. The usage of pooled sequencing offers a scalable solution for population analysis and for species level comparisons studies. This study may serve as a basis for future mitochondrial genome variability studies and representations.

# 1 First steps towards mitochondrial 2 pan-genomics: Detailed analysis of 3 *Fusarium graminearum* mitogenomes

4 Balázs Brankovics<sup>1,2,3</sup>, Tomasz Kulik<sup>4</sup>, Jakub Sawicki<sup>4</sup>, Katarzyna Bilka<sup>4</sup>,  
5 Hao Zhang<sup>5</sup>, G Sybren de Hoog<sup>2,3</sup>, Theo AJ van der Lee<sup>1</sup>, Cees Waalwijk<sup>1</sup>,  
6 and Anne D van Diepeningen<sup>1,2</sup>

7 <sup>1</sup>B.U. Biointeractions and Plant Health, Wageningen University and Research,  
8 Wageningen, Netherlands

9 <sup>2</sup>Westerdijk Fungal Biodiversity Institute, Royal Netherlands Academy of Arts and  
10 Sciences, Utrecht, Netherlands

11 <sup>3</sup>Institute for Biodiversity and Ecosystem Dynamics, University of Amsterdam,  
12 Amsterdam, Netherlands

13 <sup>4</sup>Department of Botany and Nature Protection, University of Warmia and Mazury,  
14 Olsztyn, Poland

15 <sup>5</sup>State Key Laboratory for Biology of Plant Diseases and Insect Pests, Institute of Plant  
16 Protection, Chinese Academy of Agriculture Sciences, Beijing, China P.R.

17 Corresponding author:

18 Balázs Brankovics<sup>1,2,3</sup>

19 Email address: balazs.brankovics@wur.nl

## 20 ABSTRACT

21 There is a gradual shift from representing a species' genome by a single reference genome sequence  
22 to a pan-genome representation. Pan-genomes are the abstract representations of the genomes of all  
23 the strains that are present in the population or species. In this study, we employed a pan-genomic  
24 approach to analyze the intraspecific mitochondrial genome diversity of *Fusarium graminearum*. We  
25 present an improved reference mitochondrial genome for *F. graminearum* with an intron-exon annotation  
26 that was verified using RNA-seq data. Each of the 24 studied isolates had a distinct mitochondrial  
27 sequence. Length variation in the *F. graminearum* mitogenome was found to be largely due to variation  
28 of intron regions (99.98%). The "intronless" mitogenome length was found to be quite stable and  
29 could be informative when comparing species. The coding regions showed high conservation, while  
30 the variability of intergenic regions was highest. However, the most important variable parts are the  
31 intron regions, because they contain approximately half of the variable sites, make up more than half  
32 of the mitogenome, and show presence/absence variation. Furthermore, our analyses show that the  
33 mitogenome of *F. graminearum* is recombining, as was previously shown in *F. oxysporum*, indicating that  
34 mitogenome recombination is a common phenomenon in *Fusarium*. The majority of mitochondrial introns  
35 in *F. graminearum* belongs to group I introns, which are associated with homing endonuclease genes  
36 (HEGs). Mitochondrial introns containing HE genes may spread within populations through homing,  
37 where the endonuclease recognizes and cleaves the recognition site in the target gene. After cleavage  
38 of the "host" gene, it is replaced by the gene copy containing the intron with HEG. We propose to use  
39 introns unique to a population for tracking the spread of the given population, because introns can spread  
40 through vertical inheritance, recombination as well as via horizontal transfer. We demonstrated how  
41 pooled sequencing of strains can be used for mining mitogenome data. The usage of pooled sequencing  
42 offers a scalable solution for population analysis and for species level comparisons studies. This study  
43 may serve as a basis for future mitochondrial genome variability studies and representations.

## 44 INTRODUCTION

45 One of the most ideal markers for monitoring the distribution and spread of populations is the mitochon-  
46 drial genome (Harrison, 1989). Due to its high copy number within individual cells, the mitochondrial  
47 genome is easy to access. Furthermore, it is mostly maternally inherited and it is less likely to recombine  
48 than the nuclear genome. In fungi gender is not genetically determined, and since maternal structures  
49 and meiosis require resources, the better adapted genotype is more likely to act as the maternal strain.  
50 This means that the mitochondrial genotype has the potential to be used to track the successful nuclear  
51 genotypes.

52 Mitochondrial sequences have been used for resolving phylogenetic and evolutionary relationships  
53 between fungi at all taxonomic levels (Liu et al., 2009; Avila-Adame et al., 2006; Fourie et al., 2013).  
54 In 2003, the DNA barcoding initiative started, aiming at using a single marker for taxon identification.  
55 The marker that was selected was a mitochondrial gene, cytochrome c oxidase I – COI or *cox1* (Hebert  
56 et al., 2003). In *Fusarium* however, the use of *cox1* was abandoned as a barcoding region, because of the  
57 frequent presence of introns in the gene made this region impractical for PCR amplification (Gilmore et al.,  
58 2009). Next generation sequencing (NGS) and new analysis methods have resolved this issue (Brankovics  
59 et al., 2016).

60 *Fusarium graminearum* is the major causative agent of Fusarium head blight (FHB), a devastating  
61 disease of small grain cereals. Besides reducing yield, the fungus contaminates crops with mycotoxins  
62 such as trichothecenes and zearalenone, which pose a serious threat to food and feed safety (Desjardins,  
63 2006). Population studies of *F. graminearum* showed that the populations are highly dynamic and several  
64 displacements have been reported (Gale et al., 2007; Ward et al., 2008). Monitoring these population  
65 shifts is important, as they may differ in virulence, fungicide resistance and/or mycotoxin profile (Gale  
66 et al., 2007; Zhang et al., 2012).

67 The mitochondrial genome of *F. graminearum* encodes all genes typically associated with mtDNAs  
68 of filamentous fungi: two rRNA coding genes, 14 protein coding genes and a large set of tRNA coding  
69 genes (Al-Reedy et al., 2012). In addition, a large open reading frame with unknown function (LV-  
70 uORF) was found, flanked by tRNA genes. The first comparative studies of mitochondrial genomes of  
71 *Fusarium* spp. have revealed that *F. graminearum* has a significantly larger mitogenome than *Fusarium* spp.  
72 belonging to other species complexes analyzed so far (Fourie et al., 2013; Al-Reedy et al., 2012). Intron  
73 variation within the FGSC has not yet been analyzed, but the mitogenomes of different species within the  
74 *F. fujikuroi* species complex showed diversity in intron content based on the sequences of *F. circinatum*,  
75 *F. fujikuroi* and *F. verticillioides* (Fourie et al., 2013).

76 Most mitochondrial introns found in *Fusarium* are group I introns. These introns are self-splicing  
77 ribozymes, which frequently contain homing endonuclease genes (HEGs) (Haugen et al., 2005). The  
78 combination of intron and HEG forms a mobile element that is able to invade intronless copies of the  
79 “host” gene (Haugen et al., 2005), thereby enabling horizontal spread of the mobile element through the  
80 population. This mechanism is called homing, since the homing endonuclease recognizes a target site  
81 of 15–45 bp, which makes the insertion highly sequence specific (Haugen et al., 2005). A functional  
82 homing endonuclease is needed for the homing of the intron, but the intron may be retained as long  
83 as the self-splicing function of the intron is intact. Since the mitochondrial genes are crucial for the  
84 proper functioning of the cell, if an intron loses its ability to self-splice, then the intron is lost through  
85 precise excision (Goddard and Burt, 1999). This mechanism allows an intron to spread in populations to  
86 strains that do not possess the given intron. This dispersion does not require further recombination. The  
87 mechanism does not allow one haplotype of an intron to replace another one, since the horizontal transfer  
88 is mediated only by the cleavage of an intronless copy. Hence, the replacement of one haplotype by  
89 another one can only be explained either by recombination or by loss of the original intron and insertion  
90 of the new haplotype.

91 Pan-genomes are the abstract representation of the genomes of all the strains that are present in the  
92 population. The idea of pan-genome or supra-genome comes from bacterial genomics, and originated  
93 from the distributed genome hypothesis (DGH) (Ehrlich, 2001; Tettelin et al., 2005). According to the  
94 DGH, each strain within a population/species contains a subset of contingency genes from within the  
95 supra-genome (pan-genome), i.e., the supra-genome is distributed among many individual strains (Ehrlich,  
96 2001; Ehrlich et al., 2004). Pan-genome based analysis can be used to identify conserved, variable and  
97 strain specific regions within a group of genomes. Pan-genomes can be also employed to contrast two  
98 populations or two species.

99 In order to create a pan-genome for the mitogenome of *F. graminearum*, we have to better understand  
100 the nature and dynamics of the diversity in the mitochondrial genome of this organism. To accomplish  
101 this, a reliable reference has to be established as a basis for all comparative analyses. To this end, we  
102 resequenced the reference strain of *F. graminearum*, PH-1, assembled its mitochondrial genome, improved  
103 its annotation and validated the annotation using RNA-seq. Subsequently, this reference was used to  
104 study the SNP frequencies, intron distribution and sequence variability of the different regions of the  
105 mitogenome within the species, by analyzing a total of 24 strains, which were individually sequenced,  
106 representing a wide range of hosts and geographic origins. Finally, we evaluated the efficacy of using  
107 pooled sequencing in assessing the mitogenome sequence diversity within a sample. Pooled sequencing  
108 offers the possibility of analyzing populations directly from field samples.

## 109 MATERIALS & METHODS

### 110 Strains

111 Thirteen *F. graminearum* strains were sequenced individually on the Illumina *Miseq* platform (Table 1).  
112 In addition, *F. graminearum* strain PH-1 (CBS 123657, NRRL 31084) was sequenced on the Illumina  
113 *HiSeq* platform both as a single strain and as part of a pooled set of five *F. graminearum* strains (Table 1).  
114 Besides the newly sequenced strains, the whole genome sequencing reads of ten *F. graminearum* were  
115 downloaded from the SRA database of NCBI that were produced by other research groups (Laurent et al.,  
116 2017; Wang et al., 2017). The outgroup, *F. gerlachii* strain was sequenced for an earlier publication (Kulik  
117 et al., 2016). A detailed description of the fungal strains is given in Table 1.

### 118 Sequencing

#### 119 *Illumina Miseq*

120 Whole genome libraries were prepared using the Nextera XT kit (Illumina, San Diego, CA, USA) from  
121 gDNA extracted from mycelium. The constructed libraries were sequenced on the Illumina *Miseq* platform  
122 with 250 bp paired-end read, version 2. The fungal genomes were sequenced in a multiplexed format (6-7  
123 samples per run), where an oligonucleotide index barcode was embedded within adapter sequences that  
124 were ligated to DNA fragments (Smith et al., 2010). Next, the sequence reads were de-multiplexed and  
125 filtered for low quality base calls, trimming all bases from 5' and 3' read ends with Phred scores <Q30.

#### 126 *Illumina HiSeq*

127 For *F. graminearum* strain PH-1 (CBS 123657, NRRL 31084) a random sheared shotgun library was  
128 prepared using the NEXTflex ChIP-seq Library prep kit with adaptations for low input gDNA according  
129 to manufacturer's protocol (Bioscientific). The library was loaded as (part of) one lane of an Illumina  
130 paired-end flowcell for cluster generation using a cBot. Sequencing was done on an Illumina HiSeq2000  
131 instrument using 101, 7, 101 flow cycles for forward, index and reverse reads respectively. De-multiplexing  
132 of resulting data was carried out using the Casava 1.8 software. Sequencing reads have been uploaded to  
133 the European Nucleotide Archive (ENA) with the accession number PRJEB18592.

134 The same method was applied for the pooled sequencing with the adjustment that random sheared  
135 shotgun library was prepared by using equal amounts of genomic DNA extract from all five strains  
136 (Table 1). Sequencing reads have been uploaded to the European Nucleotide Archive (ENA) with the  
137 accession number PRJEB18596.

#### 138 *Third party sequencing data*

139 Besides the sequencing data that we have generated, we also made use of sequencing data produced  
140 by other research groups that had been submitted to SRA (Sequencing Read Archive) databases. This  
141 included a dataset of SRA data of 6 strains isolated from France (PRJNA295638; Laurent et al., 2017), 3  
142 strains from China (PRJNA296400; Wang et al., 2017) and one strain from Australia (PRJNA235346;  
143 Gardiner et al., 2014). The mitochondrial genome sequences for the strains sequenced by third party  
144 are available in the Third Party Annotation Section of the DDBJ/ENA/GenBank databases under the  
145 accession numbers TPA: BK010538-BK010547

### 146 Assembly

147 GRABb was used with SPAdes assembler to reconstruct the mitogenome of the strains. GRABb (Brankovics  
148 et al., 2016) was chosen because it is a wrapper program for iterative *de novo* assembly based on a refer-  
149 ence sequence. SPAdes 3.8.1 (Bankevich et al., 2012; Nurk et al., 2013) assembler was used, since



150 it offers good insight for the user into the relationship between nodes in the assembly graph and the  
151 relationship between nodes, contigs and scaffolds. The mitochondrial genomes were assembled from  
152 NGS reads using GRABb by specifying the mitogenome sequence of PH-1 strain (HG970331) as query  
153 sequence.

154 For each individually sequenced strain it was possible to resolve the assembly to a single circular  
155 sequence. When the GRABb run finished for the strains that were pooled for sequencing, the final  
156 assembly graph was visualized using Bandage (Wick et al., 2015) and the assembly was resolved to two  
157 circular sequence variants to capture all the variation within the dataset (Supplementary Text 1). For the  
158 first circular sequence, referred to as “short”, the shorter alternative contigs were included in the path at  
159 each position where continuity was ambiguous. While for the other sequence, referred to as “long”, the  
160 longer alternatives were included. In this way, all the different sequence regions were represented at least  
161 once in the two sequences.

### 162 **Sequence annotation**

163 The initial mitogenome annotations were done using MFannot (<http://megasun.bch.umontreal.ca/cgi-bin/mfannot/mfannotInterface.pl>)  
164 and were manually improved: annotation of tRNA genes was improved using tRNAscan-SE (Pavesi  
165 et al., 1994), annotation of protein-coding genes and the *rnl* gene was corrected by aligning intron-  
166 less homologs to the genome. Intron encoded proteins were identified using NCBI’s ORF Finder  
167 (<http://www.ncbi.nlm.nih.gov/gorf/gorf.html>) and annotated using InterPro (Mitchell  
168 et al., 2015) and CD-Search (Marchler-Bauer and Bryant, 2004). The annotated mitochondrial genome  
169 sequences are available under the following GenBank accession numbers: BK010538-BK010547,  
170 KP966550-KP966561, KR011238 and MH412632.

### 172 **Read mapping and SNP discovery**

173 The mitogenome of *F. graminearum* strain PH-1 and the two mitogenome sequences obtained from the  
174 assembly of the pooled dataset were used as reference sequences for the read mapping and SNP discovery.  
175 The read mapping was done using *aln* and *sampe* subcommands of the Burrows-Wheeler Alignment tool  
176 (BWA-0.7.12-r1034) (Li and Durbin, 2009). SNP calling was done using SAMtools mpileup (1.3.1) with  
177 *-g* and *-f* flag and BCFtools call (1.3.1) with *-mv* flag (Li et al., 2009).

### 178 **Coverage analysis**

179 Coverage of different regions was estimated by, first, mapping reads of the pooled dataset to the reference  
180 sequence using the *sampe* subcommand of the Burrows-Wheeler Alignment tool (BWA-0.7.12-r1034) (Li  
181 and Durbin, 2009). Then, read coverage was calculated using the *genomecov* command of bedtools  
182 v2.26.0. The following single copy nuclear protein coding genes were used to represent single copy  
183 nuclear regions:  $\gamma$ -actin (*act*),  $\beta$ -tubulin II (*tub2*), calmodulin (*cal*), 60S ribosomal protein L10 (*rpl10a*),  
184 the second largest subunit of DNA-dependent RNA polymerase II (*rpb2*), translation elongation factor  
185  $1\alpha$  (*tef1a*), translation elongation factor 3 (*tef3*) and topoisomerase I (*top1*). The reference sequences  
186 were extracted from the genome of PH-1 (4 chromosomes: HG970332, HG970333, HG970334, and  
187 HG970335). The nuclear mitochondrial DNA segment (NUMT) used for coverage comparison was  
188 identified during the assembly of the pooled data (see Supplementary Text 1).

### 189 **Intron validation**

190 The RNA-seq data for *F. graminearum* PH-1 was downloaded from NCBI’s SRA database, accession  
191 number PRJNA239711 (Zhao et al., 2014). Read mapping was done by *subjunc* command of the Subread  
192 aligner (Liao et al., 2013). Intron positions were identified from the CIGAR string of the SAM file  
193 produced by the aligner.

### 194 **Linear model**

R was used for linear model analysis to test whether the intron variation is the main reason of mitochondrial  
genome length variation within the species. The linear model was the following:

$$y = x + c$$

where  $y$  was the total length of the mitochondrial genome,  $x$  was the length of the intron sequences and  
 $c$  was the  $y$ -intercept (average intronless length of the mitochondrial genomes). The  $R^2$  value obtained

from linear model analysis specifies what percentage of the variation of the dependent value (mitogenome length) is explained by the variation in the independent value (intron length).

$$R^2 = 1 - \frac{SS_{residual}}{SS_{total}}$$

195 Residual sums of squares ( $SS_{residual}$ ) and total sums of squares ( $SS_{total}$ ) were calculated using the *deviance*  
196 function of R.

### 197 **Comparative sequence analysis**

198 The nucleotide sequences were aligned using MUSCLE (Edgar, 2004b,a). Sequence variability of given  
199 regions was calculated by aligning the sequences. Then the number of characters with multiple character  
200 states was calculated and divided by the total number of characters in the alignment. This step was  
201 done using *fasta\_variability* from the *fasta\_tools* package ([https://github.com/b-brankovics/fasta\\_tools](https://github.com/b-brankovics/fasta_tools)).  
202

### 203 **Phylogenetic analysis**

204 The most appropriate substitution evolution model was determined using jModelTest 2 (Darriba et al.,  
205 2012) for each of the regions analyzed. Phylogenetic trees were calculated using RAxML version  
206 8.2.4 (Stamatakis, 2014). Two measures of clade support were used in this study: i) bootstrap (BS) values  
207 calculated by 1000 bootstrap runs using RAxML and ii) Bayesian posterior probability (BPP). In order to  
208 obtain BPP values, phylogenetic reconstruction has been conducted using MrBayes 3.2.5 (Ronquist et al.,  
209 2012). The MCMC algorithm was run for 4,000,000 generations with four incrementally-heated chains,  
210 starting from random trees and sampling one out every 1000 generations. Burn-in was set to relative  
211 burn-in of 0.25. The generated tree-space was used to calculate the BPP.

### 212 **Detecting the presence of recombination**

213 The intergenic regions were analyzed using the  $\Phi_w$ -test implemented in SplitsTree (Bruen et al., 2006) to  
214 detect whether there is recombination in the mitochondrial genome.

## 215 **RESULTS**

### 216 **Mitochondrial genome of *F. graminearum***

217 The mitochondrial genomes of all 24 strains sequenced individually were assembled into single circular  
218 contigs. The re-sequencing of the mitochondrial genome of *F. graminearum* strain PH-1 revealed two  
219 SNPs compared to the most recent published mitogenome assembly (HG970331.1) of the strain that  
220 was based on next generation sequencing reads (King et al., 2015). The correction of these SNPs was  
221 supported by the fact that all the other strains contained the same two SNPs obtained in the new assembly  
222 of PH-1. The newly assembled mitochondrial genome of the PH-1 strain as well as the other mitochondrial  
223 genomes were annotated. The mitochondrial genomes of all strains contained the same set of genes in  
224 the same order and orientation (Fig. 1). To test whether the intron-exon models were predicted correctly,  
225 RNA-seq reads were mapped against the mitogenome of *F. graminearum* strain PH-1. The results of the  
226 read mapping supported all of the predicted intron-exon boundaries.

### 227 **Mitogenome variability in *F. graminearum***

228 The mitogenomes of *F. graminearum* strains analyzed showed variation in size, ranging from 93,560 bp  
229 to 101,424 bp (Table 2). To test whether intron variation is the main reason of mitochondrial genome  
230 length variation within the species, linear model analysis was used. The linear model that assumed that  
231 mitochondrial length variation is due only to variation of the length of intron regions explained 99.98% of  
232 intraspecific length variation observed in the data, showing that intron variation is the main reason behind  
233 intraspecific mitochondrial genome length variation. The standard deviation of the mitogenome length  
234 was 1818 bp, which is 1.87% of the average mitochondrial genome length.

235 The coding regions (tRNA, rRNA and conserved protein coding genes) showed low levels of variation  
236 both within *F. graminearum* (0.02%) and when compared to *F. gerlachii* (0.02%). In addition, none of the  
237 SNPs found in protein coding regions caused amino acid substitution.

238 The large ORF with unknown function (LV-uORF) located in the large variable region of the mi-  
239 togenome contained five SNPs within *F. graminearum* and the sequence in the *F. gerlachii* strain was

240 identical to the most frequent haplotype within *F. graminearum*. All five SNPs resulted in amino acid  
241 substitution in the putative peptide sequences. The variability of the conserved protein coding regions was  
242 0.02%, while the variability of the LV-uORF region was 0.09% within *F. graminearum*. The difference  
243 in variability was even more striking on the protein sequence level, where the conserved protein genes  
244 showed no variation, while the LV-uORF showed 0.26% variability.

245 The variability of the intergenic regions was 1.63% and 2.30% for intraspecies and interspecies,  
246 respectively. The overall sequence variability of intron sequences was 0.68% and 0.71% for intraspecies  
247 and interspecies, respectively. Although the variability of intron regions was significantly less than that of  
248 intergenic regions, both regions contained approximately equal numbers of variable sites (Table 3) due to  
249 the large length difference between the two regions. The intron regions were the most variable part of  
250 the mitochondrial genomes, because approximately half of the variable sites was located in introns, and  
251 introns were the only regions showing presence/absence variation within *F. graminearum*.

252 Interestingly, strains CBS 128539 and CBS 138561 had identical intergenic sequences, while strains  
253 CBS 104.09 and CBS 119800 (isolated 81 years apart) had identical intron sequences. However, all the  
254 strains had a unique mitochondrial genome sequence.

### 255 **Intron patterns and phylogeny**

256 A total of 39 intron sites were found in the individually sequenced dataset (Supplementary Table ). Out  
257 of the 39 introns, 32 were present in all strains and 21 of these showed no variation at the intraspecies  
258 level and 14 at the interspecies level. The introns that showed presence/absence variation within the  
259 dataset were cob-i159, cob-i201, cox1-i1287, cox2-i228, cox2-i318, cox2-i552 and nad2-i1632 (Fig. 2  
260 and Supplementary Table ). The intron names contain the gene name where they are located and the  
261 coding nucleotide position of the host gene after which they were inserted.

262 It was not possible to group the strains based on their intron patterns (presence/absence for each  
263 intron) without allowing for multiple gain or loss of introns (Supplementary Table ). This could be the  
264 result of recombination of the mitochondrial genome or the horizontal transfer of introns. Recombination  
265 would affect intergenic regions, while the horizontal transfer of the intron by homing would not affect the  
266 intergenic regions. Recombination of the intergenic regions was well supported ( $p = 2.26 * 10^{-6}$ ) by the  
267  $\Phi_w$ -test.

### 268 **Strategies to analyze pooled mitochondrial NGS data**

269 Two approaches were used to explore the mitogenome variability in the pooled dataset: i) assembling the  
270 reads *de novo* and ii) mapping the reads to a reference sequence.

#### 271 ***De novo assembly approach***

272 The assembly resulted in a graph that contained five ambiguous sites that represented four inser-  
273 tion/deletion variations (three intron presence/absence variation, cob-i201, cox1-i1287, cox2-i318, and  
274 a large insertion inside the cob-i490 intron) in the dataset, and one site (located in nad4L-i239) where  
275 two different alleles were found in the strain set (Supplementary Text 1). These polymorphic sites were  
276 too far apart to establish linkage between them, so two alternative assemblies were extracted from the  
277 assembly graph: one with the shorter allele at all of the positions and one with the longer allele at all of  
278 the positions (Supplementary Text 1). The assembly method did not reveal SNP variations, only intron  
279 presence/absence variations and one replacement variation.

#### 280 ***Mapping approach***

281 To assess the influence of the reference sequence on the mapping and SNP calling results, both of the  
282 sequences obtained from the assembly approach of the pooled dataset were used as references, beside  
283 the curated mitogenome of the PH-1 strain. Besides giving an insight into the influence of the reference  
284 sequence to the downstream analysis, this also makes it possible to detect variation within intron sequences  
285 that may be absent in some of the reference sequences.

286 The lowest coverage detected for a single nucleotide allele was 21% of the reads that mapped to the  
287 given position. This is close to the expected value (20%) for an allele present in a single strain in a pool  
288 of five strains. This result shows that the method was sensitive enough to detect a SNP present in a single  
289 strain. Furthermore, the results of all three analyses identified the same polymorphic sites. This means  
290 that the choice of reference sequence did not influence the SNP detection results.

291 The three runs of read mapping and SNP calling revealed a total of fifteen SNPs (Table 4). The allele  
292 ratios were identical even when the reference sequence used for the mapping was different, with one

293 exception: position 90636. At this position both PH-1 and the pooled assembly analysis showed 70%  
294 for the nucleotide present in the given reference and 30% for the alternative, despite the fact that the two  
295 references had different nucleotides at the given location (Table 4). Examination of the alignment of the  
296 reference sequences revealed that the sequence difference was not only a single nucleotide polymorphism  
297 at position 90636, but there was a 8 bp long indel at position 90627-90634. This nearby indel influenced  
298 the mapping of reads containing the allele differing from the reference sequence. This was the reason  
299 why the SNP calling skewed in favor of the reference allele in both mappings.

### 300 **Coverage analysis**

301 Coverage values were calculated for different genomic regions in order to determine whether coverage  
302 cutoffs could be used to differentiate between mitochondrial and nuclear regions. The coverage of single  
303 copy nuclear regions that were present in all of the pooled strains was 290x. The coverage of the nuclear  
304 mitochondrial DNA segment (NUMT) sequence was 230x, which suggests that it was present in four of  
305 the five pooled strains. The coverage of mitogenome regions that were present in all strains was 4000x.  
306 While, the coverage of singleton mitochondrial regions, present only in a single strain, was 475x. The  
307 coverage gap was sufficiently high between shared single copy nuclear regions (290x) and singleton  
308 mitochondrial sequence (475x) to allow clear differentiation between them.

## 309 **DISCUSSION**

310 Comparative genomics analyses are traditionally reference (Laurent et al., 2017) or pairwise based (Fourie  
311 et al., 2013). Reference based methods are efficient at identifying regions that are present in the reference,  
312 but absent in other individuals, or detecting smaller variations, like SNPs. This method does not identify  
313 regions that are absent from the single reference, while these regions might be valuable for clustering the  
314 non-reference individuals. Pairwise comparison is able to identify unique regions for both individuals;  
315 however, it is difficult to scale to a larger sample size, because every individual has to be compared to  
316 every other individual, then the results have to be brought to the same scale.

317 To take full advantage of next generation sequencing data, a paradigm shift is needed: from focusing  
318 on a single reference genome to using a pan-genome, that is, a representation of all genomic content in  
319 a certain population, species or phylogenetic clade (Computational Pan-Genomics Consortium, 2018).  
320 In this study, we used an *ad hoc* pan-genomic analysis of the mitochondrial genomes of *Fusarium*  
321 *graminearum*. The reason for using an *ad hoc* approach is that pan-genomics is still a young field of  
322 research, and as such, there are no clear standards developed yet for analysis, for files or for data sharing.  
323 The goal of the analysis was to understand the nature and the dynamics of mitogenome variability, then  
324 to identify the implications of these results for mitogenome based population studies or track & trace  
325 implementations. The results of this study can be utilized for the development of suitable data structures  
326 and file formats for capturing the variability of mitochondrial pan-genomes.

327 In this study, we improved the mitochondrial genome reference for *F. graminearum* strain PH-1,  
328 which is recognized as the reference strain of this species for genomic studies (Al-Reedy et al., 2012;  
329 King et al., 2015; Cuomo et al., 2007). The first mitochondrial genome sequence was produced using  
330 Sanger sequencing and primer walking by Al-Reedy et al. (2012). The assembly was improved by King  
331 et al. (2015) using NGS reads. This assembly corrected 15 SNPs and 30 indels in the sequence. Here, we  
332 present a new assembly, which corrected 2 more SNPs, complete with a detailed annotation. The introns  
333 that were predicted during the annotation process were all verified by RNA-seq data.

334 The mitochondrial genomes of *F. graminearum* and *F. gerlachii* contained the same genes and ORFs  
335 in the same orientation. The coding sequences showed high levels of conservations, and all SNPs found in  
336 protein coding genes were synonymous substitution. The genetic variation in the mitochondrial genome  
337 could be classified in two groups: small sequence variations (SNPs and short indels) and intron gain  
338 and loss. Although, variations resulting from SNPs and short indels were twice as frequent in intergenic  
339 regions as in intron regions, about half of the variable sites was located in intron regions. The second  
340 type of variation, the presence/absence of introns, accounted for 99.98% of the length variation between  
341 the mitochondrial genomes. In conclusion, the majority of the sequence variation within the species  
342 was related to intron regions: either SNPs and short indels or the presence/absence of complete introns.  
343 Thus, in mitogenome comparative analysis or pan-genomic studies, special attention should be given to  
344 accurately capturing the intron variation, since it is the most informative fraction of the mitogenome.

345 An alternative way to sequencing strains individually is sequencing them in a pool. The pooled

346 sequencing approach is more cost efficient than sequencing the strains separately. The data produced by  
347 pooled sequencing of strains from a given population could be viewed as the pan-genomic sequencing  
348 reads of that population. In this study, we have demonstrated how sequencing data from pools of strains  
349 can be mined for mitochondrial genome variation. Sequencing in pools has already been used to discover  
350 rare alleles of nuclear loci (Raineri et al., 2012). This method can be used for finding rare alleles, but  
351 it also allows a scalable solution for analyzing complete populations. So far, the application of pooled  
352 sequencing data has been used for SNP discovery in nuclear loci from samples (Raineri et al., 2012).  
353 However, analyzing mitochondrial genome data of fungi possesses some additional challenges. We have  
354 demonstrated that besides SNPs, intron presence/absence variation is a major element of the mitogenome  
355 variation. To assess what kind of approach can detect intron presence/absence variation and SNP variation,  
356 we analyzed the data using a *de novo* assembly approach followed by a read mapping and SNP-calling  
357 approach. The results show that the assembly approach is able to identify sequence differences affecting  
358 sequence regions longer than individual sequencing reads, such as insertions and deletions of intron  
359 sequence or long polymorphic sequences, while it is unable to identify SNPs or short indels. Read  
360 mapping and SNP calling analysis has to be performed to identify SNPs. This method in turn is unable to  
361 identify sequence differences affecting longer sequence regions. For optimal results, a sequential approach  
362 is needed for analyzing pooled samples: first, an assembly step to identify introns or larger indels absent  
363 from the reference genome, then using both the reference and the newly identified extra regions for read  
364 mapping and SNP-calling.

365 The drawbacks of pooled data are that short indel variation might be missed and linkage between  
366 markers is lost when using short read sequencing technologies, although linkage information is not crucial  
367 when comparing pan-genomes with each other. Furthermore, pooling large amount of strains could mean  
368 the loss of the coverage gap between mitochondrial copies and nuclear copies, which makes NGS analysis  
369 of mitochondrial genomes more advantageous to PCR methods. This means that nuclear mitochondrial  
370 sequences (NUMTs) might affect the results. With sufficient caution the effects of NUMTs can be  
371 minimized, since they can be identified in the assembly step. In the assembly step, NUMTs may appear as  
372 separate contigs, as in our example, or as new paths similar to introns with the significant difference that  
373 intron segments are joined to the rest of the mitochondrial assembly on both termini, while the flanking  
374 nuclear regions of NUMTs would only be joined on one of the termini of the segment. Despite these  
375 concerns, the benefits of pooled sequencing of large numbers of strains offers a scalable solution for  
376 population or species level comparisons. After a reference sequence is established, each population or  
377 species could then be represented by pools of multiple strains.

378 Most of the introns in *F. graminearum* are group I introns, and contain homing endonuclease genes  
379 (HEGs). Group I introns harboring a functional HEG can spread in a population through homing. Homing  
380 is facilitated by the homing endonuclease that cleaves the target gene at a 15-45 bp recognition site. The  
381 resulting double strand break stimulates homologous recombination based DNA repair. Since all copies of  
382 the mitochondrial genome that contain the recognition site are susceptible to the homing endonuclease, the  
383 only viable template for homologous repair is a genome that contains a copy of the intron. The insertion  
384 of the intron into the recognition site modifies the sequence, and it will no longer be recognized by the  
385 homing endonuclease.

386 The mitochondrial genome of *F. graminearum* shows evidence of recombination. We recently showed  
387 that mitochondrial recombination does also happen in the *F. oxysporum* species complex (Brankovics et al.,  
388 2017). Recombination of the mitochondrial genome in *Fusarium* appears to be a common phenomenon,  
389 since both *F. oxysporum* and *F. graminearum* show signs of mitochondrial recombinations, despite the  
390 fact that *F. oxysporum* is an asexual fungus with a putative parasexual cycle, while *F. graminearum* is  
391 a homothallic species that has an active sexual cycle (Yun et al., 2000). Due to the recombination of  
392 the mitogenome, it cannot directly be used as a marker to track successful nuclear genotypes as was  
393 proposed. However, based on the spreading mechanism of introns, introns could be used for track and  
394 trace implementations. The intron sequences spread through clonal & sexual reproductions, and through  
395 horizontal transfer. Due to the effect of the homing endonuclease, all offspring of a sexual cross would be  
396 tagged by all the introns that are specific to either parent. The appearance of a new intron in a population  
397 could be a sign of migration or gene flow.

398 The annotation of strain CBS 119173 revealed a putative nested intron in *cox1-i906*. All other strains  
399 contain a haplotype that is 1006 bp long, while this strain contains a haplotype that is 2084 bp long. The  
400 sequence comparison indicates that the additional 1078 bp region is an intron that was integrated inside

401 the homing endonuclease of the acceptor intron. This putative intron contains an additional HEG, but  
 402 the annotation pipeline did not identify the sequence as an intron. This indicates that introns and intron  
 403 encoded genes themselves are susceptible for intron invasions. The question is whether the invading  
 404 intron has to retain its self-splicing function or the “host” (or acceptor) intron can splice the complete  
 405 nested construct with its own self-splicing activity.

406 The intron regions contain most of the variation within *F. graminearum* and population specific introns  
 407 promise to be valuable markers for tracking.

## 408 CONCLUSIONS

409 We have improved the reference mitochondrial genome [reference](#) sequence for *F. graminearum*. Intraspe-  
 410 cific mitochondrial genome length variations are mainly due to intron presence/absence variation, thus  
 411 using “intronless” length—subtracting the length of the intron regions from the total mitogenome length—  
 412 could be a valuable information when comparing species. Mitogenomes are also subject to recombination  
 413 in both *F. graminearum* and in *F. oxysporum*, indicating that it is a common phenomenon in *Fusarium*.  
 414 We proposed that introns unique to a single population could be used to track the spread of the given  
 415 population, because introns can spread through vertical inheritance, recombination and horizontal transfer.  
 416 We also demonstrated how pooled sequencing of strains can be used for the mitogenome. The usage of  
 417 pooled sequencing offers a scalable solution for population analysis and for species level comparisons  
 418 [studies](#). The results of this study represent an important step towards establishing pan-genomics for  
 419 mitochondrial genomes.

## 420 REFERENCES

- 421 Al-Reedy, R. M., Malireddy, R., Dillman, C. B., and Kennell, J. C. (2012). Comparative analysis of  
 422 *Fusarium* mitochondrial genomes reveals a highly variable region that encodes an exceptionally large  
 423 open reading frame. *Fungal Genetics and Biology*, 49(1):2–14.
- 424 Avila-Adame, C., Gómez-Alpizar, L., Zismann, V., Jones, K. M., Buell, C. R., and Ristaino, J. B.  
 425 (2006). Mitochondrial genome sequences and molecular evolution of the Irish potato famine pathogen,  
 426 *Phytophthora infestans*. *Current Genetics*, 49(1):39–46.
- 427 Bankevich, A., Nurk, S., Antipov, D., Gurevich, A. A., Dvorkin, M., Kulikov, A. S., Lesin, V. M.,  
 428 Nikolenko, S. I., Pham, S., Prjibelski, A. D., Pyshkin, A. V., Sirotkin, A. V., Vyahhi, N., Tesler, G.,  
 429 Alekseyev, M. A., and Pevzner, P. A. (2012). SPAdes: A new genome assembly algorithm and its  
 430 applications to single-cell sequencing. *Journal of Computational Biology*, 19(5):455–477.
- 431 Brankovics, B., van Dam, P., Rep, M., de Hoog, G. S., van der Lee, T. A. J., Waalwijk, C., and van  
 432 Diepeningen, A. D. (2017). Mitochondrial genomes reveal recombination in the presumed asexual  
 433 *Fusarium oxysporum* species complex. *BMC Genomics*, 18(1):735.
- 434 Brankovics, B., Zhang, H., van Diepeningen, A. D., van der Lee, T. A. J., Waalwijk, C., and de Hoog,  
 435 G. S. (2016). GRAB: Selective assembly of genomic regions, a new niche for genomic research. *PLoS  
 436 Comput Biol.*, 12(6):e1004753.
- 437 Bruen, T. C., Philippe, H., and Bryant, D. (2006). A simple and robust statistical test for detecting the  
 438 presence of recombination. *Genetics*, 172(4):2665–2681.
- 439 Computational Pan-Genomics Consortium (2018). Computational pan-genomics: status, promises and  
 440 challenges. *Briefings in bioinformatics*, 19(1):118–135.
- 441 Cuomo, C. A., Güldener, U., Xu, J.-R., Trail, F., Turgeon, B. G., Di Pietro, A., Walton, J. D., Ma, L.-J.,  
 442 Baker, S. E., Rep, M., Adam, G., Antoniw, J., Baldwin, T., Calvo, S., Chang, Y.-L., Decaprio, D., Gale,  
 443 L. R., Gnerre, S., Goswami, R. S., Hammond-Kosack, K., Harris, L. J., Hilburn, K., Kennell, J. C.,  
 444 Kroken, S., Magnuson, J. K., Mannhaupt, G., Mauceli, E., Mewes, H.-W., Mitterbauer, R., Muehlbauer,  
 445 G., Münsterkötter, M., Nelson, D., O’Donnell, K., Ouellet, T., Qi, W., Quesneville, H., Roncero, M.  
 446 I. G., Seong, K.-Y., Tetko, I. V., Urban, M., Waalwijk, C., Ward, T. J., Yao, J., Birren, B. W., Kistler,  
 447 H. C., O’donnell, K., Ouellet, T., Qi, W., Quesneville, H., Roncero, M. I. G., Seong, K.-Y., Tetko, I. V.,  
 448 Urban, M., Waalwijk, C., Ward, T. J., Yao, J., Birren, B. W., and Kistler, H. C. (2007). The *Fusarium  
 449 graminearum* genome reveals a link between localized polymorphism and pathogen specialization.  
 450 *Science*, 317(5843):1400–1402.
- 451 Darriba, D., Taboada, G. L., Doallo, R., and Posada, D. (2012). jModelTest 2: more models, new  
 452 heuristics and parallel computing. *Nature Methods*, 9(8):772–772.

- 453 Desjardins, A. E. (2006). *Fusarium Mycotoxins: Chemistry, Genetics, and Biology*. The American  
454 Phytopathological Society, St. Paul, MN, USA.
- 455 Edgar, R. C. (2004a). MUSCLE: a multiple sequence alignment method with reduced time and space  
456 complexity. *BMC Bioinformatics*, 19(5):113.
- 457 Edgar, R. C. (2004b). MUSCLE: multiple sequence alignment with high accuracy and high throughput.  
458 *Nucleic Acids Res.*, 32(5):1792–1797.
- 459 Ehrlich, G. D. (2001). The biofilm and distributed genome paradigms provide a new theoretical structure  
460 for understanding chronic bacterial infections. In *Interscience Conference on Antimicrobials Agents  
461 and Chemotherapy (ICAAC)*, Chicago, IL, USA.
- 462 Ehrlich, G. D., Hu, F. Z., and Post, J. C. (2004). Role for biofilms in infectious disease. In Ghannoum, M.  
463 and O'Toole, G. A., editors, *Microbial biofilms*, chapter 18, pages 332–358. ASM Press, Washington,  
464 DC.
- 465 Fourie, G., van der Merwe, N. A., Wingfield, B. D., Bogale, M., Tudzynski, B., Wingfield, M. J., and  
466 Steenkamp, E. T. (2013). Evidence for inter-specific recombination among the mitochondrial genomes  
467 of *Fusarium* species in the *Gibberella fujikuroi* complex. *BMC genomics*, 14(1):605.
- 468 Gale, L. R., Ward, T. J., Balmas, V., and Kistler, H. C. (2007). Population subdivision of *Fusarium*  
469 *graminearum sensu stricto* in the Upper Midwestern United States. *Phytopathology*, 97(11):1434–1439.
- 470 Gardiner, D. M., Stiller, J., and Kazan, K. (2014). Genome sequence of *Fusarium graminearum* isolate  
471 CS3005. *Genome Announcements*, 2(2):e00227–14.
- 472 Gilmore, S. R., Gräfenhan, T., Louis-Seize, G., and Seifert, K. A. (2009). Multiple copies of cytochrome  
473 oxidase 1 in species of the fungal genus *Fusarium*. *Molecular Ecology Resources*, 9(Suppl. 1):90–98.
- 474 Goddard, M. R. and Burt, A. (1999). Recurrent invasion and extinction of a selfish gene. *Proceedings of  
475 the National Academy of Sciences of the United States of America*, 96(24):13880–13885.
- 476 Harrison, R. G. (1989). Animal mitochondrial DNA as a genetic marker in population and evolutionary  
477 biology. *Trends in Ecology and Evolution*, 4(1):6–11.
- 478 Haugen, P., Simon, D. M., and Bhattacharya, D. (2005). The natural history of group I introns. *Trends in  
479 Genetics*, 21(2):111–119.
- 480 Hebert, P. D. N., Cywinska, A., Ball, S. L., and DeWaard, J. R. (2003). Biological identifications through  
481 DNA barcodes. *Proceedings of the Royal Society B: Biological Sciences*, 270(1512):313–321.
- 482 King, R., Urban, M., Hammond-Kosack, M. C. U., Hassani-Pak, K., and Hammond-Kosack, K. E. (2015).  
483 The completed genome sequence of the pathogenic ascomycete fungus *Fusarium graminearum*. *BMC  
484 Genomics*, 16:544.
- 485 Kulik, T., Brankovics, B., Sawicki, J., and van Diepeningen, A. D. (2016). The complete mitogenome of  
486 *Fusarium gerlachii*. *Mitochondrial DNA. Part A, DNA mapping, sequencing, and analysis*, 27(3):1895–  
487 6.
- 488 Laurent, B., Moïnard, M., Spataro, C., Ponts, N., Barreau, C., and Foulongne-Oriol, M. (2017). Landscape  
489 of genomic diversity and host adaptation in *Fusarium graminearum*. *BMC genomics*, 18(203):1–19.
- 490 Li, H. and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform.  
491 *Bioinformatics*, 25(14):1754–1760.
- 492 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., and Durbin,  
493 R. (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, 25(16):2078–2079.
- 494 Liao, Y., Smyth, G. K., and Shi, W. (2013). The Subread aligner: Fast, accurate and scalable read mapping  
495 by seed-and-vote. *Nucleic Acids Research*, 41(10):e108.
- 496 Liu, Y., Steenkamp, E. T., Brinkmann, H., Forget, L., Philippe, H., and Lang, B. F. (2009). Phylogenomic  
497 analyses predict sistergroup relationship of nucleariids and Fungi and paraphyly of zygomycetes with  
498 significant support. *BMC evolutionary biology*, 9:272.
- 499 Marchler-Bauer, A. and Bryant, S. H. (2004). CD-Search: Protein domain annotations on the fly. *Nucleic  
500 Acids Research*, 32(Web Server issue):327–331.
- 501 Mitchell, A., Chang, H.-Y., Daugherty, L., Fraser, M., Hunter, S., Lopez, R., McAnulla, C., McMenamin,  
502 C., Nuka, G., Pesseat, S., Sangrador-Vegas, A., Scheremetjew, M., Rato, C., Yong, S.-Y., Bateman, A.,  
503 Punta, M., Attwood, T. K., Sigrist, C. J. A., Redaschi, N., Rivoire, C., Xenarios, I., Kahn, D., Guyot,  
504 D., Bork, P., Letunic, I., Gough, J., Oates, M., Haft, D., Huang, H., Natale, D. A., Wu, C. H., Orengo,  
505 C., Sillitoe, I., Mi, H., Thomas, P. D., and Finn, R. D. (2015). The InterPro protein families database:  
506 The classification resource after 15 years. *Nucleic Acids Research*, 43(Database issue):D213–D221.
- 507 Nurk, S., Bankevich, A., Antipov, D., Gurevich, A. A., Korobeynikov, A., Lapidus, A., Prjibelski, A. D.,

- 508 Pyshkin, A., Sirotkin, A., Sirotkin, Y., Stepanauskas, R., Clingenpeel, S. R., Woyke, T., McLean,  
509 J. S., Lasken, R., Tesler, G., Alekseyev, M. A., and Pevzner, P. A. (2013). Assembling single-cell  
510 genomes and mini-metagenomes from chimeric MDA products. *Journal of Computational Biology*,  
511 20(10):714–37.
- 512 Pavesi, A., Conterio, F., Bolchi, A., Dieci, G., and Ottonello, S. (1994). Identification of new eukaryotic  
513 tRNA genes in genomic DNA databases by a multistep weight matrix analysis of transcriptional control  
514 regions. *Nucleic Acids Res.*, 22(7):1247–1256.
- 515 Raineri, E., Ferretti, L., Esteve-Codina, A., Nevado, B., Heath, S., and Pérez-Enciso, M. (2012). SNP  
516 calling by sequencing pooled samples. *BMC Bioinformatics*, 13:239.
- 517 Ronquist, F., Teslenko, M., van der Mark, P., Ayres, D. L., Darling, A., Höhna, S., Larget, B., Liu, L.,  
518 Suchard, M. A., and Huelsenbeck, J. P. (2012). MrBayes 3.2: Efficient bayesian phylogenetic inference  
519 and model choice across a large model space. *Syst. Biol.*, 61(3):539–542.
- 520 Smith, A. M., Heisler, L. E., St.Onge, R. P., Farias-Hesson, E., Wallace, I. M., Bodeau, J., Harris,  
521 A. N., Perry, K. M., Giaever, G., Pourmand, N., and Nislow, C. (2010). Highly-multiplexed barcode  
522 sequencing: an efficient method for parallel analysis of pooled samples. *Nucleic Acids Research*,  
523 38(13):e142.
- 524 Stamatakis, A. (2014). RAxML version 8: A tool for phylogenetic analysis and post-analysis of large  
525 phylogenies. *Bioinformatics*, 30(9):1312–1313.
- 526 Tettelin, H., Masignani, V., Cieslewicz, M. J., Donati, C., Medini, D., Ward, N. L., Angiuoli, S. V.,  
527 Crabtree, J., Jones, A. L., Durkin, A. S., DeBoy, R. T., Davidsen, T. M., Mora, M., Scarselli, M.,  
528 Margarit y Ros, I., Peterson, J. D., Hauser, C. R., Sundaram, J. P., Nelson, W. C., Madupu, R., Brinkac,  
529 L. M., Dodson, R. J., Rosovitz, M. J., Sullivan, S. A., Daugherty, S. C., Haft, D. H., Selengut, J., Gwinn,  
530 M. L., Zhou, L., Zafar, N., Khouri, H., Radune, D., Dimitrov, G., Watkins, K., O'Connor, K. J. B.,  
531 Smith, S., Utterback, T. R., White, O., Rubens, C. E., Grandi, G., Madoff, L. C., Kasper, D. L., Telford,  
532 J. L., Wessels, M. R., Rappuoli, R., and Fraser, C. M. (2005). Genome analysis of multiple pathogenic  
533 isolates of *Streptococcus agalactiae*: Implications for the microbial "pan-genome". *Proceedings of the*  
534 *National Academy of Sciences*, 102(39):13950–13955.
- 535 Wang, Q., Jiang, C., Wang, C., Chen, C., Xu, J.-R., and Liu, H. (2017). Characterization of the two-speed  
536 subgenomes of *Fusarium graminearum* reveals the fast-speed subgenome specialized for adaption and  
537 infection. *Frontiers in Plant Science*, 8:140.
- 538 Ward, T. J., Clear, R. M., Rooney, A. P., O'Donnell, K., Gaba, D., Patrick, S., Starkey, D. E., Gilbert,  
539 J., Geiser, D. M., and Nowicki, T. W. (2008). An adaptive evolutionary shift in *Fusarium* head blight  
540 pathogen populations is driving the rapid spread of more toxigenic *Fusarium graminearum* in North  
541 America. *Fungal Genetics and Biology*, 45:473–484.
- 542 Wick, R. R., Schultz, M. B., Zobel, J., and Holt, K. E. (2015). Bandage: Interactive visualization of *de*  
543 *novo* genome assemblies. *Bioinformatics*, 31(20):3350–3352.
- 544 Yun, S. H., Arie, T., Kaneko, I., Yoder, O. C., and Turgeon, B. G. (2000). Molecular organization  
545 of mating type loci in heterothallic, homothallic, and asexual *Gibberella/Fusarium* species. *Fungal*  
546 *genetics and biology*, 31(1):7–20.
- 547 Zhang, H., van der Lee, T. A. J., Waalwijk, C., Chen, W., Xu, J., Xu, J., Zhang, Y., and Feng, J. (2012).  
548 Population analysis of the species complex from wheat in China show a shift to *Fusarium graminearum*  
549 more aggressive isolates. *PLoS One*, 7(2):e31722.
- 550 Zhao, C., Waalwijk, C., de Wit, P., Tang, D., and van der Lee, T. (2014). Relocation of genes generates  
551 non-conserved chromosomal segments in *Fusarium graminearum* that show distinct and co-regulated  
552 gene expression patterns. *BMC Genomics*, 15(1):191.



553 **ADDITIONAL FILES**554 **Supplementary Table 1**

555 **Intron locations, lengths and haplotypes within standard mitochondrial genes of *Fusarium graminearum* and *F. gerlachii* strains.** The different haplotypes are displayed with different background color  
 556 per intron site. Haplotypes that have identical length are differentiated from each other by using ', ^ or \*,  
 557 corresponding to the number of SNPs differentiating the haplotypes (haplotypes 1159' and 1159'' differ  
 558 by 2 SNPs, while 1159 and 1159''' differ by 3 SNPs).

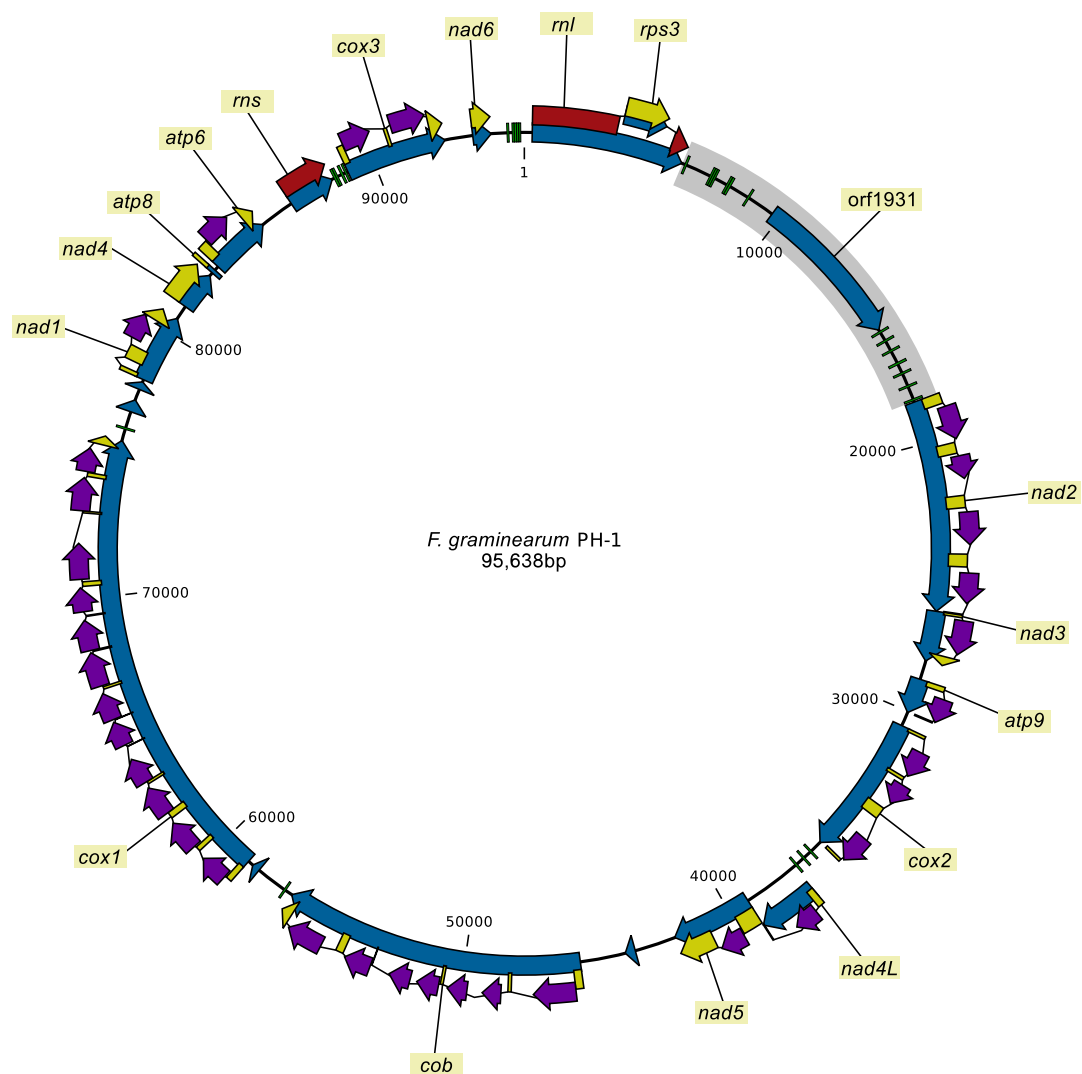
560 **Supplementary Text 1**

561 Assembling the pooled data set.

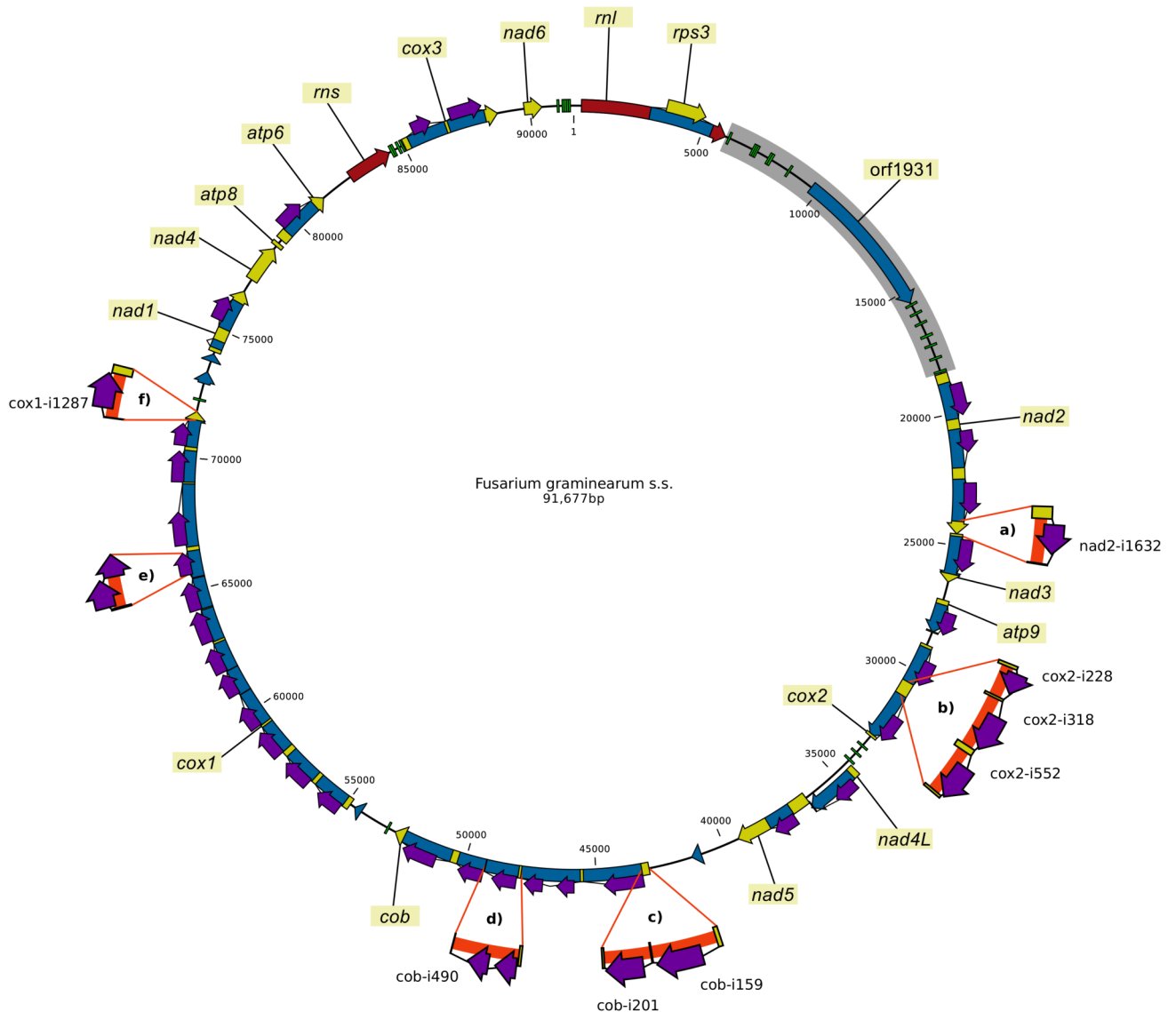
562 **TABLES AND FIGURES**

**Table 1. List of *Fusarium* strains analysed in this study**

Species	Strain	Origin	Host	Year of isolation	Sequenced individually or in a pool
<i>F. graminearum</i>	CBS123657 (PH-1) NRRL31084	USA	maize	1996	both
<i>F. graminearum</i>	CBS119173	USA	wheat head	2005	individually
<i>F. graminearum</i>	CBS139513	Argentina	barley	2011	individually
<i>F. graminearum</i>	CBS139514	Argentina	barley	2010	individually
<i>F. graminearum</i>	CBS119799	South Africa	wheat kernel	1987	individually
<i>F. graminearum</i>	CBS119800	South Africa	maize	1990	individually
<i>F. graminearum</i>	CBS110263	Iran	maize	1968	individually
<i>F. graminearum</i>	CBS123688	Sweden	oats	unknown	individually
<i>F. graminearum</i>	CBS128539	Belgium	wheat kernel	2007	individually
<i>F. graminearum</i>	CBS138561	Poland	wheat kernel	2010	individually
<i>F. graminearum</i>	CBS138562	Poland	wheat kernel	2010	individually
<i>F. graminearum</i>	CBS138563	Poland	wheat kernel	2003	individually
<i>F. graminearum</i>	CBS104.09	unknown	unknown	1909	individually
<i>F. graminearum</i>	CBS185.32	unknown	maize	1932	individually
<i>F. graminearum</i>	CS3005	Australia	barley	2001	individually
<i>F. graminearum</i>	HN9-1	China	wheat	2002	individually
<i>F. graminearum</i>	HN-Z6	China	wheat	2012	individually
<i>F. graminearum</i>	INRA-156	France	wheat	2001	individually
<i>F. graminearum</i>	INRA-159	France	wheat	2001	individually
<i>F. graminearum</i>	INRA-164	France	wheat	2002	individually
<i>F. graminearum</i>	INRA-171	France	wheat	2001	individually
<i>F. graminearum</i>	INRA-181	France	wheat	2002	individually
<i>F. graminearum</i>	INRA-195	France	wheat	2002	individually
<i>F. graminearum</i>	YL-1	China	wheat	2012	individually
<i>F. graminearum</i>	bfb0999.1	China	barley	2005	pooled
<i>F. graminearum</i>	68D2	Netherlands	wheat	2001	pooled
<i>F. graminearum</i>	CHG013	China	maize	2005	pooled
<i>F. graminearum</i>	CHG157	China	barley	2005	pooled
<i>F. gerlachii</i>	CBS123666	USA	wheat head	2000	individually



**Figure 1.** The mitogenome of *F. graminearum* strain PH-1. Green blocks: tRNA coding genes, blue arrows: genes or ORFs (no labels added for short ORFs), yellow arrows: protein coding sequences, red arrows: rDNA coding sequence, purple arrows: intron encoded homing endonuclease genes, gray box: the large variable (LV) region with orf1931 (LV-uORF).



**Figure 2. Pan-genomic representation of the presence/absence variation of introns in the mitochondrial genomes of the 24 *F. graminearum* strains.**

In the figure, the thick orange lines highlight intron sequences in the alternative sequences. (SNPs and short indels are not indicated.) **a)** The insertion of nad2-i1632; **b)** the insertion of cox2-i228, cox2-i318 and cox2-i552; **c)** the insertion of cob-i159 and cob-i201; **d)** longer variant of cob-i490; **e)** intron insertion in the HEG located in cox1-i906; and **f)** the insertion of cox1-i1287.

**Table 2. Mitochondrial genome variation of the *Fusarium graminearum* strains**

Strain ID	GenBank accession numbers	Size (bp)	Introns	Intronic (bp)	Core (bp)
CBS123657 (PH-1)	MH412632	95638	34	49429	46209
CBS185.32	KP966550	96300	34	50120	46180
CBS110263	KP966551	97364	35	51165	46199
CBS119173	KP966552	100342	37	54130	46212
CBS119799	KP966553	96005	35	49919	46086
CBS119800	KP966554	97462	35	51280	46182
CBS123688	KP966555	95035	34	48837	46198
CBS128539	KP966556	96134	35	49996	46138
CBS138561	KP966557	95034	34	48837	46197
CBS138562	KP966558	99062	36	52865	46197
CBS138563	KP966559	99068	36	52865	46203
CBS139514	KP966560	96167	35	49980	46187
CBS139513	KP966561	95041	34	48837	46204
CBS104.09	KR011238	97460	35	51280	46180
CS3005	BK010538	93560	33	47381	46179
HN9-1	BK010539	96307	35	51567	44740
HN-Z6	BK010540	97767	34	50120	47647
INRA-156	BK010541	101424	37	55243	46181
INRA-159	BK010542	96199	35	49996	46203
INRA-164	BK010543	99678	37	53476	46202
INRA-171	BK010544	96199	35	49996	46203
INRA-181	BK010545	96187	35	49996	46191
INRA-195	BK010546	97358	35	51165	46193
YL-1	BK010547	97996	36	51777	46219

**Core** stands for the total mitogenome length minus the length of the intron regions.

**Table 3. Distribution of variation in the intron and intergenic regions within and between species**

	Intraspecies			Interspecies		
	Length (bp)	Variable positions	Variation frequency	Length (bp)	Variable positions	Variation frequency
Coding	21572	4	0.02%	21572	5	0.02%
intron	59091	399	0.68%	59091	419	0.71%
Intergenic	18982	310	1.63%	18982	436	2.30%

**Table 4. List of single nucleotide polymorphisms identified in the pooled dataset of *Fusarium graminearum* strains.** Positions are aligned positions between the PH-1 reference sequence and the pooled sequences (“short” and “long”). “Reference” refers to the nucleotide found in the given reference sequence used for mapping, while “Alternative” refers to the nucleotide suggested by the mapped reads. Position 90636 shows unusual ratios: in both mappings the reference nucleotide (C or A) has a frequency of 70% and the alternative nucleotide has 30%. This is due to an adjacent indel that affects the mapping results.

Position	PH-1		Pooled	
	Reference	Alternative	Reference	Alternative
2337	A (0.77)	G (0.23)	A (0.77)	G (0.23)
6288	C (0.41)	A (0.59)	A (0.61)	C (0.39)
6355	T (0.42)	C (0.58)	C (0.60)	T (0.40)
13540	C (0.78)	A (0.22)	C (0.78)	A (0.22)
37126	C (0.75)	T (0.25)	C (0.75)	T (0.25)
37773	A (0.75)	G (0.25)	A (0.75)	G (0.25)
44773	A (0.62)	G (0.38)	A (0.62)	G (0.38)
64776	G (0.53)	A (0.47)	G (0.53)	A (0.47)
70827	A (0.62)	G (0.38)	A (0.62)	G (0.38)
89194	G (0.57)	A (0.43)	G (0.57)	A (0.43)
90636	C (0.70)	A (0.30)	A (0.70)	C (0.30)
95918	A (0.43)	C (0.57)	C (0.59)	A (0.41)
99784	A (0.40)	G (0.60)	G (0.62)	A (0.38)
100362	C (0.42)	A (0.58)	A (0.59)	C (0.41)
100538	G (0.42)	A (0.58)	A (0.61)	G (0.39)