

ORIGINAL RESEARCH

The insertion of a mitochondrial selfish element into the nuclear genome and its consequences

Julien Y. Dutheil^{1,2,3}  | Karin Münch²  | Klaas Schotanus^{2,4}  |
Eva H. Stukenbrock^{1,2,4}  | Regine Kahmann² 

¹Max Planck Institute for Evolutionary Biology, Plön, Germany

²Max Planck Institute for Terrestrial Microbiology, Marburg, Germany

³Institute of Evolutionary Sciences, CNRS – University of Montpellier – IRD – EPHE, Montpellier, France

⁴Christian Albrechts University of Kiel, Kiel, Germany

Correspondence

Julien Y. Dutheil, Max Planck Institute for Evolutionary Biology, August-Thienemann-Str. 2, 24306 Plön, Germany.
Email: dutheil@evolbio.mpg.de

Present address

Klaas Schotanus, Department of Molecular Genetics and Microbiology (MGM), Duke University Medical Center, Durham, NC, USA

Funding information

Max-Planck-Gesellschaft

Abstract

Homing endonucleases (HE) are enzymes capable of cutting DNA at highly specific target sequences, the repair of the generated double-strand break resulting in the insertion of the HE-encoding gene (“homing” mechanism). HEs are present in all three domains of life and viruses; in eukaryotes, they are mostly found in the genomes of mitochondria and chloroplasts, as well as nuclear ribosomal RNAs. We here report the case of a HE that accidentally integrated into a telomeric region of the nuclear genome of the fungal maize pathogen *Ustilago maydis*. We show that the gene has a mitochondrial origin, but its original copy is absent from the *U. maydis* mitochondrial genome, suggesting a subsequent loss or a horizontal transfer from a different species. The telomeric HE underwent mutations in its active site and lost its original start codon. A potential other start codon was retained downstream, but we did not detect any significant transcription of the newly created open reading frame, suggesting that the inserted gene is not functional. Besides, the insertion site is located in a putative RecQ helicase gene, truncating the C-terminal domain of the protein. The truncated helicase is expressed during infection of the host, together with other homologous telomeric helicases. This unusual mutational event altered two genes: The integrated HE gene subsequently lost its homing activity, while its insertion created a truncated version of an existing gene, possibly altering its function. As the insertion is absent in other field isolates, suggesting that it is recent, the *U. maydis* 521 reference strain offers a snapshot of this singular mutational event.

KEYWORDS

gene birth, gene transfer, homing endonuclease, intron, mitochondrion

1 | INTRODUCTION

The elucidation of the mechanisms at the origin of genetic variation is a longstanding goal of molecular evolutionary biology.

Mutation accumulation experiments—together with comparative analysis of sequence data—are instrumental in studying the processes shaping genetic diversity at the molecular level (Eyre-Walker & Keightley, 2007; Kondrashov & Kondrashov, 2010).

This article has been peer-reviewed and recommended by PCI Evolutionary Biology (<https://doi.org/10.24072/pci.evolbiol.100101>).

This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2020 The Authors. *Ecology and Evolution* published by John Wiley & Sons Ltd

They revealed that the spectrum of mutations ranges from single nucleotide substitutions to large scale chromosomal rearrangements, and encompasses insertions, deletions, inversions, and duplication of genetic material of variable length (Lynch et al., 2008). Mutation events may result from intrinsic factors such as replication errors and repair of DNA damage. In some cases, however, mutations can be caused or favored by extrinsic factors, such as mutagenic environmental conditions or parasitic genome entities like viruses or selfish mobile elements. Such particular sequences, able to replicate and invade the host genome, may have various effects including inserting long stretches of DNA that do not encode any organismic function, but also disrupting, copying, and moving parts of the genome sequence. These selfish element-mediated mutations can significantly contribute to the evolution of their host: First, the invasion of mobile elements creates “junk” DNA that can significantly increase the genome size (Lynch, 2007), and some of this material can be ultimately domesticated and acquire a new function, beneficial to the host (Kaessmann, 2010; Volff, 2006). Second, the genome dynamics resulting from the activity of mobile elements can generate novelty by gene duplication (Dutheil et al., 2016; Ohta, 2000) or serve as a mechanism of parasexuality and compensate for the reduced diversity in the absence of sexual reproduction (Dong, Raffaele, & Kamoun, 2015; Möller & Stukenbrock, 2017). Finally, control mechanisms (such as repeat-induced point mutations in fungi (Gladyshev, 2017)) may also incidentally affect genetic diversity (Grandaubert et al., 2014).

Intron-borne homing endonuclease genes (HEGs) constitute a class of selfish elements whose impact on genome evolution is less well documented. They encode a protein able to recognize a particular genomic DNA sequence and cut it (homing endonuclease, HE). The resulting double-strand break is subsequently repaired by recombination using the homologous sequence containing the HEG itself as a template, resulting in its insertion in the target location (Stoddard, 2005). As the recognized sequence is typically large, its occurrence is rare and the insertion typically happens at a homologous position. In this process, a *heg*⁺ element containing the endonuclease gene converts a *heg*⁻ allele (devoid of HEG but harboring the recognition sequence) to *heg*⁺, a mobility mechanism referred to as *homing* (Dujon et al., 1989). After the insertion, the host cell is homozygous *heg*⁺, and the HEG segregates at a higher frequency than the Mendelian rate (Goddard & Burt, 1999). The open reading frame of the HEG is typically associated with a sequence capable of self-splicing, either at the RNA (group-I introns) or protein (inteins) level, avoiding disruption of functionality when inserted in a protein-coding gene (Chevalier & Stoddard, 2001; Stoddard, 2005). The dynamic of HEGs has been well described and involves three stages: (a) conversion from *heg*⁻ to *heg*⁺ by homing activity, (b) degeneration of the HEG leading to the loss of homing activity, but still protecting against a new insertion because the target is altered by the insertion event, and (c) loss of the HEG leading to the restoration of the *heg*⁻ allele (Barzel, Obolski, Gogarten, Kupiec, & Hadany, 2011; Gogarten & Hilario, 2006). This cycle leads to recurrent gains and losses of HEG at a given genomic position, and ultimately to the loss of the

HEG at the population level unless new genes invade from other locations or by horizontal gene transfer (Gogarten & Hilario, 2006).

Homing endonuclease genes are found in all domains of life and in the genomes of organelles, mitochondria, and chloroplasts (Belfort & Roberts, 1997; Lambowitz & Belfort, 1993; Stoddard, 2005). In several fungi, HEGs are residents of mitochondria. Here, we study the molecular evolution of a HEG from the fungus *Ustilago maydis*, which serves as a model for the elucidation of (a) fundamental biological processes like cell polarity, morphogenesis, organellar targeting, and (b) the mechanisms allowing biotrophic fungi to colonize plants and cause disease (Ast, Stiebler, Freitag, & Böcker, 2013; Djamei & Kahmann, 2012; Steinberg & Perez-Martin, 2008). *U. maydis* is the most well-studied representative of smut fungi, a large group of plant pathogens, because of the ease by which it can be manipulated both genetically and through reverse genetics approaches (Vollmeister et al., 2012). Besides, its compact, fully annotated genome comprises only 20.5 Mb and is mostly devoid of repetitive DNA (Kämper et al., 2006). The genome sequences of several related species, *Sporisorium reilianum*, *S. scitamineum*, and *U. hordei* causing head smut in corn, smut whip in sugarcane and covered smut in barley, respectively, provide a powerful resource for comparative studies (Dutheil et al., 2016; Laurie et al., 2012; Schirawski et al., 2010). Here, we report the case of a mitochondrial HEG that integrated into the nuclear genome of *U. maydis*. This singular mutation event created two new genes: First, the original endonuclease activity of the integrated HEG was inactivated by a deletion in the active site, leading to a frameshift and a new open reading frame containing the DNA-binding domain of the HEG (Derbyshire, Kowalski, Dansereau, Hauer, & Belfort, 1997). Second, the integration of the HEG occurred within another protein-coding gene, leading to its truncation.

2 | MATERIALS AND METHODS

2.1 | Analysis of codon usage and GC content

Ustilago maydis gene models (genome version 2.0) were retrieved from the MIPS database (Mewes et al., 2011). Mitochondrial genes were extracted from the *U. maydis* full mitochondrial genome (Genbank accession number: NC_008368.1). Within-group correspondence analysis of synonymous codon usage was performed using the *ade4* package for R, following the procedure described in (Charif, Thioulouse, Lobry, & Perrière, 2005). The proportion of G and C nucleotides was computed along with the first 10 kb of *U. maydis* chromosome 9, using 300 bp windows slid by 1 bp.

2.2 | Strains, growth conditions, and virulence assays

The *Escherichia coli* strains DH5 α (Bethesda Research Laboratories) and TOP10 (Life Technologies) were used for the cloning and amplification of plasmids. *U. maydis* strains 518 and 521 are the parents

of FB1 and FB2 (Banuett & Herskowitz, 1989). SG200 is a haploid solopathogenic strain derived from FB1 (Kämper et al., 2006). 10–1 is an uncharacterized haploid *U. maydis* strain isolated in the United States and kindly provided by G. May. I2, O2, P2, S5, and T6 are haploid *U. maydis* strains collected in different parts of Mexico (Valverde, Vandemark, Martínez, & Paredes-López, 2000). The haploid *S. reilianum* strains SRZ1 and SRZ2 as well as the solopathogenic strain JS161 derived from SRZ1 have been described (Schirawski et al., 2010). Deletion mutants were generated by gene replacement using a PCR-based approach and verified by Southern analysis (Kämper, 2004).

pRS426 Δ um11064 + 11065 is a pRS426-derived plasmid containing the *UMAG_11064/UMAG_11065* double deletion construct which consists of a hygromycin resistance cassette flanked by the left border of the *UMAG_11064* and right border of the *UMAG_11065* gene. The left border of *UMAG_11064* and the right border of *UMAG_11065* were PCR amplified from SG200 gDNA with primers um11064_lb_fw/um11064_lb_rv and um11065_rb_fw/um11065_rb_rv (Table S7). The hygromycin resistance cassette was obtained from *Sfi*I digested pHwtFRT (Khrunyk, Münch, Schipper, Lupas, & Kahmann, 2010). The pRS426 *Eco*RI/*Xho*I backbone, both borders and the resistance cassette, was assembled using yeast drag and drop cloning (Christianson, Sikorski, Dante, Shero, & Hieter, 1992). The fragment containing the deletion cassette was amplified from this plasmid using primers um11064_lb_fw and um11065_rb_rv (Table S7), transformed into SG200, and transformants carrying a deletion of *UMAG_11064* and *UMAG_11065* were identified by southern analysis (Figure S3).

Ustilago maydis strains were grown at 28°C in liquid YEPSL medium (0.4% yeast extract, 0.4% peptone, 2% sucrose) or on PD solid medium (2.4% Potato Dextrose broth, 2% agar). Stress assays were performed as described in (Krombach, Reissmann, Kreibich, Bochen, & Kahmann, 2018). Transformation and selection of *U. maydis* transformants followed published procedures (Kämper et al., 2006). To assess virulence, seven-day-old maize seedlings of the maize variety Early Golden Bantam (Urban Farmer) were syringe-infected. At least three independent infections were carried out, and disease symptoms were scored according to Kämper et al. (Kämper et al., 2006). Consistency of replicates was tested using a chi-squared test, and *p*-values were computed using 1,000,000 permutations. As no significant difference between replicates was observed (*p*-value = .347 for the wildtype and *p*-value = .829 for the deletion strain), observation was pooled between all replicates for each strain before being compared.

2.3 | Blast searches and gene alignment

We performed BlastN and BlastP (Altschul, Gish, Miller, Myers, & Lipman, 1990) searches using the (translated) sequence of *UMAG_11064* as a query using NCBI online blast tools. The nonredundant nucleotide and protein sequence databases were selected for BlastN and BlastP, respectively. Results were further processed

with scripts using the NCBI XML module from BioPython (Cock et al., 2009). The Macse codon aligner (Ranwez, Harispe, Delsuc, & Douzery, 2011) was used in order to infer the position of putative frameshifts in the upstream region of *UMAG_11064*. The alignment was depicted using the Boxshade software and was further manually annotated. The sequences of *U. maydis* *cox1* intron 6, as well as *S. reilianum* *cox1* introns 1 and 2, were used as query and searched against the protein nonredundant database using NCBI BlastX, excluding environmental samples and model sequences. The *cox1* genes from *U. maydis* and *S. reilianum* were aligned, and pairwise similarity was computed in nonoverlapping 100 bp windows. The gene structure, synteny, and local pairwise similarity were depicted using the genoPlotR package for R (Guy, Kultima, & Andersson, 2010).

2.4 | Phylogeny estimation, estimation of dN/dS ratios, and tests of positive selection

The nucleotide sequence of *UMAG_11064*, the first intron of the *cox1* gene of *S. reilianum*, and the eight nonredundant, most similar matches from BlastP (Table S2) were aligned using the Macse codon aligner (Ranwez et al., 2011) together with the unannotated but similar nucleotide sequences from *S. scitamineum*, *U. bromivora*, *T. indica*, and *T. walkiri*, using the NCBI codon translation table 4 “mitochondrial mold”. Columns in the alignment were manually selected to discard ambiguously aligned regions, and a phylogeny was inferred using PhyML (Guindon et al., 2010) with a general time-reversible (GTR) model of nucleotide evolution and a 4-classes discrete gamma distribution of rates. The tree topology was inferred using the “best of nearest-neighbor-interchange (NNI) and subtree-pruning-regrafting (SPR)” option, and 100 nonparametric bootstrap replicates were obtained. The final tree was rooted using the midpoint method. Analyses were performed using the Seaview software (Gouy, Guindon, & Gascuel, 2010). For the positive selection analysis, the *S. scitamineum* and *U. bromivora* sequences were discarded as they contained multiple frameshifts. A phylogenetic tree was estimated using PhyML from the remaining species after translation using a Le and Gascuel model of protein evolution (Le & Gascuel, 2008), and other parameters as for the nucleotide model. Nodes with bootstrap support values lower than 65% were collapsed. A branch model of codon evolution was fitted on the alignment and the inferred phylogenetic tree using PAML 4.9d (Yang, 2007), keeping selected positions that may contain missing data. The F3X4 codon frequency model was selected, and one dN/dS ratio was estimated per branch. A branch-site model (Zhang, Nielsen, & Yang, 2005) was fitted by specifying the branch leading to the *UMAG_11064* gene as the “foreground” group, putatively evolving under positive selection. Test for selection was performed as suggested in the PAML manual, comparing to a model where the omega2 parameter is fixed to a value of 1. A similar test was conducted after excluding the two *Tilletia* sequences from the “background” branches, as they were found to have each a branch with dN/dS > 1.

2.5 | Amplification of the *UMAG_11064* regions in several *U. maydis* strains

Amplification of DNA fragments via polymerase chain reaction (PCR) was done using the Phusion High Fidelity DNA_Polymerase (Thermo Fisher Scientific). The PCR reactions were set up in a 20 μ l reaction volume using DNA templates indicated in the respective experiments and buffer recommended by the manufacturer containing a final concentration of 3% DMSO. The PCR programs used are represented by the following scheme: Initial denaturation – [denaturation – annealing – elongation] \times number cycles – final elongation. *UMAG_11072* was amplified with primers *um11072_ORF_fw* \times *um11072_ORF_rv* using 98°C/3 m – [98°C/10 s – 65°C/30 s – 72°C/45 s] \times 30 cycles – 72°C/10 m. *UMAG_11064* was amplified with primers *um11064_ORF_fw* \times *um11064_ORF_rv* using 98°C/3 m – [98°C/10 s – 65°C/30 s – 72°C/45 s] \times 30 cycles – 72°C/10 m. The *cox1* exons 1 + 2 were amplified with primers *cox1_ex1_rv* \times *cox1_ex2_fw* using 98°C/3 m – [98°C/10 s – 63°C/30 s – 72°C/90 s] \times 33 cycles – 72°C/10 m. *cox1* exon 7 was amplified with primers *cox1_ex7_fw* \times *cox1_ex7_rv* using 98°C/3 m – [98°C/10 s – 67°C/30 s – 72°C/60 s] \times 30 cycles – 72°C/10 m. Parts of the genomic region containing *UMAG_11064*, *UMAG_11065*, and *UMAG_11066* were amplified with primer pairs *um11064_fw1* \times *um11064_rv1*, *um11064_fw1* \times *um11064_rv2*; and *um11064_fw2* \times *um11064_rv2* using 98°C/3 m – [98°C/10 s – 65°C/30 s – 72°C/150 s] \times 32 cycles – 72°C/10 m. The list of all primer sequences is provided in Table S7. PCR results are shown in Figures S1 and S2.

2.6 | History of the *UMAG_11065* family

The sequence of the *UMAG_11065* protein was used as a query for a search against several smut fungi (*U. maydis*, *U. hordei*, *S. reilianum*, *S. scitamineum*, *Melanopsichum pennsylvanicum*, and *Pseudozyma flocculosa*) complete proteome using BlastP (Altschul et al., 1990). The search finds 17 hits within the *U. maydis* genome with an *E*-value below 0.0001, as well as two genes in *S. scitamineum* (*SPSC_04622* and *SPSC_05783*) and two genes in *P. flocculosa* (*PFL1_06135* and *PFL1_02192*). Using NCBI BlastP, we found several sequences from *Fusarium oxysporum* with high similarity. We selected the sequence *FOXG_04692* as a representative and added it to the data set. The Guidance web server with the GUIDANCE2 algorithm (Sela, Ashkenazy, Katoh, & Pupko, 2015) was then used to align the protein sequences and assess the quality of the resulting alignment. Default options from the server were kept, selecting the MAFFT aligner (Katoh, Misawa, Kuma, & Miyata, 2002). Several sequences appeared to be of shallow alignment quality and were discarded. The remaining sequences were realigned using the same protocol. Four iterations were performed until the final alignment had a quality good enough for phylogenetic inference. The final alignment contained 14 sequences and had a global score of 0.79. These 14 alignable sequences contained 13 *U. maydis* sequences (including *UMAG_11065*), and the *Fusarium oxysporum* gene other sequences

from smut genomes were too divergent to be unambiguously aligned. Using Guidance, we further masked columns in the alignment with a score below 0.93 (a maximum of one position out of 14 in the column was allowed to be uncertain).

A phylogenetic analysis was conducted using the program Seaview 4 (Gouy et al., 2010). First, a site selection was performed in order to filter regions with too many gaps, leaving 506 sites. Second, a phylogenetic tree was built using PhyML within Seaview (Guindon et al., 2010) (Le and Gascuel protein substitution model (Le & Gascuel, 2008) with a four-classes discretized gamma distribution of rates, the best tree of nearest-neighbor interchange (NNI) and subtree pruning and regrafting (SPR) topological searches was kept). Support values were computed using the approximate likelihood ratio test (aLRT) method (Anisimova & Gascuel, 2006).

A test for positive selection was conducted using a combination of branch and branch-site models using PAML (Yang, 2007). The final GUIDANCE alignment was used, realigned using the Macse codon aligner (Ranwez et al., 2011), and ambiguously aligned sites and shorter sequences were manually filtered. The final alignment contained the following sequences: *UMAG_03394*, *UMAG_11065*, *UMAG_04486*, *UMAG_06506*, *UMAG_04094*, *UMAG_10585*, *UMAG_06474*, *UMAG_10980*, *UMAG_05977*, *FOXG_04692*. We used the PhyML software with the same options as described above to reconstruct a phylogenetic tree with this subset of sequences. The branch toward the *UMAG_11065* gene was used as a foreground group in the branch-site model.

2.7 | Gene expression

RNASeq normalized expression counts for the *UMAG_11064* and *UMAG_11065*, as well as of neighboring genes and paralogs elsewhere in the genome, were extracted from the Gene Expression Omnibus data set GSE103876 (Lanver et al., 2018). Gene clustering based on expression profiles was conducted using a hierarchical clustering with an average linkage on a Canberra distance, suitable for expression counts, as implemented in the “dist” and “hclust” functions in R (R Core Team, 2018). The resulting clustering tree was converted to a distance matrix and compared to the inferred phylogeny of the genes using a Mantel permutation test, as implemented in the “ape” package for R (Paradis, Claude, & Strimmer, 2004). Differences in expression between time points were assessed by fitting the linear model “expression ~ time * gene”, testing the effect of time while controlling for interaction with the “gene” variable. Residuals were normalized using a Box-Cox transform as implemented in the MASS package for R. Tukey's post hoc comparisons were conducted on the resulting model, allowing for a 5% false discovery rate.

3 | RESULTS

We report the analysis of the nuclear gene *UMAG_11064* from the smut fungus *U. maydis*, which was identified as an outlier in a

whole-genome analysis of codon usage. We first provide evidence that the gene is a former HEG and then reconstruct the molecular events that led to its insertion in the nuclear genome using comparative sequence analysis. Finally, we assess the phenotypic impact of the insertion event.

3.1 | The *UMAG_11064* nuclear gene has a mitochondrial codon usage

We studied the synonymous codon usage in protein-coding genes of the smut fungus *U. maydis*, using within-group correspondence analysis. As opposed to other methods, within-group correspondence analysis allows the comparison of codon usage while adequately taking into account confounding factors such as variation in amino-acid usage (Perrière & Thioulouse, 2002). We report a distinct synonymous codon usage for nuclear genes and mitochondrial genes (Figure 1a), with the notable exception of the nuclear gene *UMAG_11064*, which displays a typical mitochondrial codon usage. The *UMAG_11064* gene is located in the telomeric

region of chromosome 9, with no further downstream annotated gene (Figure 1b). It displays a low GC content of 30%, which contrasts with the GC content of the flanking regions (50%) and the rather homogeneous composition of the genome sequence of *U. maydis* as a whole. It is, however, in the compositional range of the mitochondrial genome (Figure 1b). Altogether, the synonymous codon usage and GC content of *UMAG_11064* suggest a mitochondrial origin.

In order to confirm the chromosomal location of *UMAG_11064*, we amplified and sequenced three regions encompassing the gene using primers within the *UMAG_11064* gene and primers in adjacent chromosomal genes upstream and downstream of *UMAG_11064* (Figure S1). The sequences of the amplified segments were in full agreement with the genome sequence of *U. maydis* (Kämper et al., 2006), thereby ruling out possible assembly artifacts in this region. As both the GC content and synonymous codon usage of *UMAG_11064* are indistinguishable from the ones of mitochondrial genes and have not moved toward the nuclear equilibrium, the transfer of the gene to its nuclear position is likely to have occurred relatively recently.

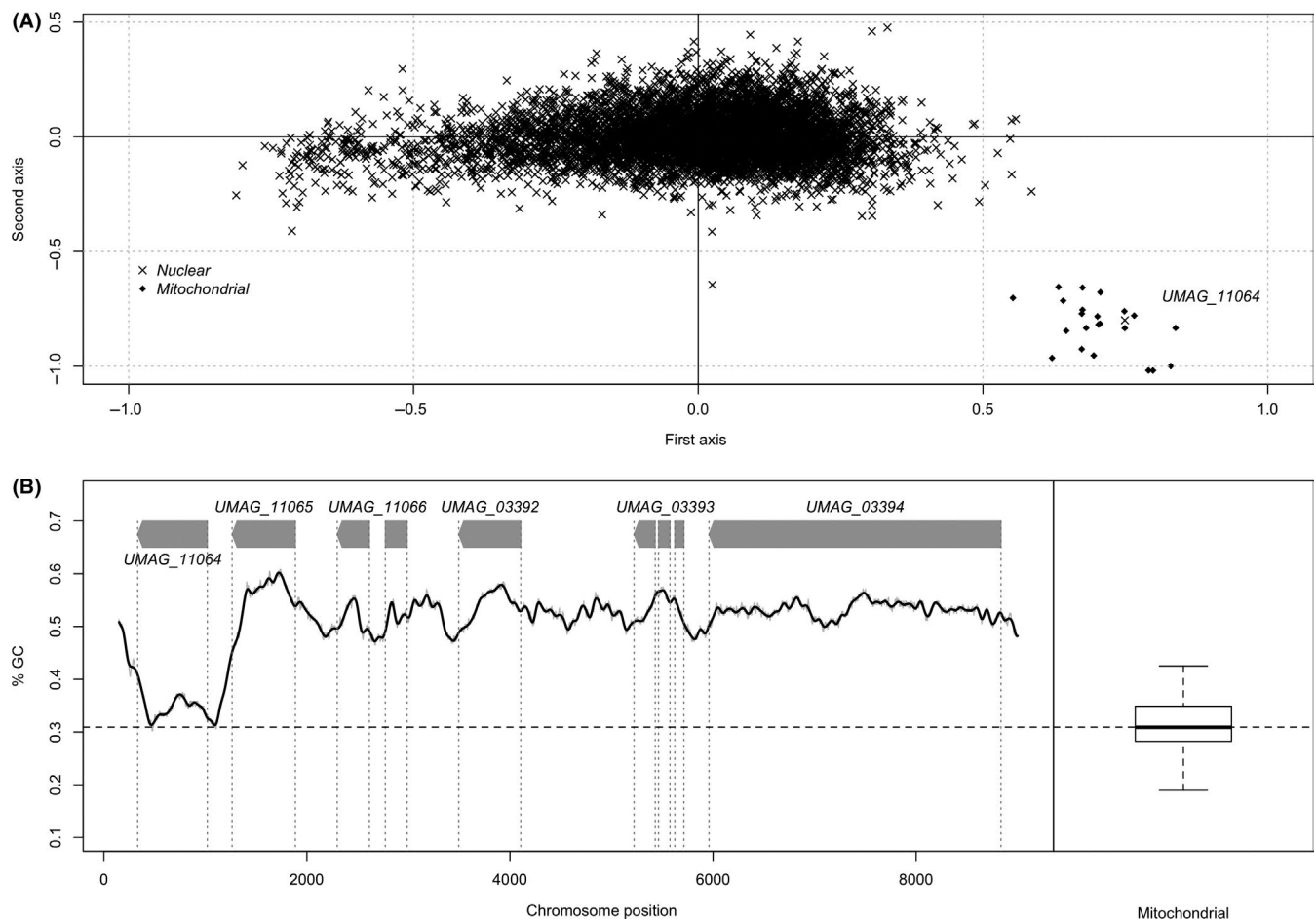


FIGURE 1 Identification of the *UMAG_11064* gene. (a) Within-group correspondence analysis of *Ustilago maydis* codon usage. Each gene is represented according to its coordinates along the first two principal factors. The genomic origin of each gene is indicated by a cross for nuclear genes and a dot for mitochondrial genes. (b) Genomic context of the gene *UMAG_11064*. GC content in 300 bp windows sliding by 1 bp, and distribution of GC content in 300 bp windows of mitochondrial genome of *U. maydis*. The dash line represents the median of the distribution

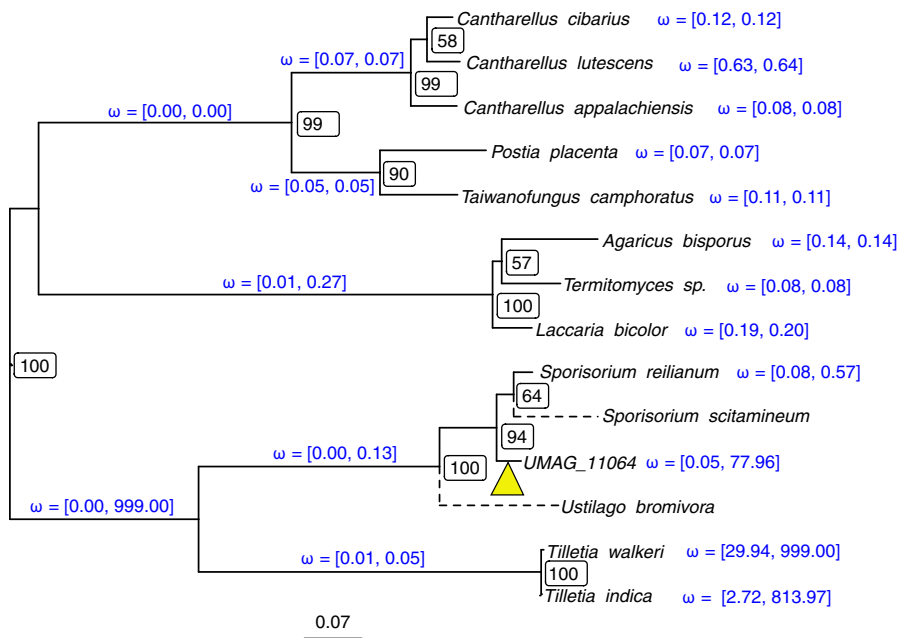


FIGURE 2 Phylogeny of *UMAG_11064* and its homologous sequences. Maximum likelihood tree inferred from nucleotide sequences. Node labels show bootstrap support values. Discontinuous branches indicate that the corresponding sequence is pseudogenized, with multiple frameshifts and insertions/deletions. Branch annotations show the minimum and maximum dN/dS ratio (ω) estimated from 10 independent runs of the codeml program. The yellow triangle indicates the supposed branch where the frameshift within the active domain of the ancestral HE occurred (see Figure 3)

3.2 | The *UMAG_11064* gene contains parts of a former GIY-YIG homing endonuclease

To gain insight into the nature of the *UMAG_11064* gene, its predicted nucleotide sequence was searched against the NCBI nonredundant nucleotide sequence database. Surprisingly, the sequence of *UMAG_11064* has no match in the mitochondrial genome of *U. maydis* itself (GenBank entry NC_008368.1), but high similarity matches were found in the mitochondrial genome of three other smut fungi (Table S1): *S. reilianum* (87% nucleotide identity), *S. scitamineum* (79%), and *U. bromivora* (76%). Two other very similar sequences were found in the mitochondrial genome of two other smut fungi, *Tilletia indica* and *Tilletia walkeri* (69% nucleotide identity), as well as in mitochondrial genomes from other basidiomycetes (e.g., *Laccaria bicolor*, 72%) and ascomycetes (e.g., *Leptosphaeria maculans*, 69%, see Table S1). The protein sequence of *UMAG_11064* shows high similarity with fungal HEGs, in particular of the so-called GIY-YIG family (Table S2) (Stoddard, 2005). The closest fully annotated protein sequence matching *UMAG_11064* corresponds to the GIY-YIG HEG located in intron 1 of the *cox1* gene of *Agaricus bisporus* (I-AbIII-P, 54% nucleotide identity). The amino-acid sequence of *UMAG_11064* matches the N-terminal part of this protein containing the DNA-binding domain of the HE (Derbyshire et al., 1997).

We performed a codon alignment of the *UMAG_11064* gene together with the most similar sequences identified by Blast, using the Macse codon aligner to infer sequence alignment in the presence of frameshifts (Ranwez et al., 2011). The sequences from *S. scitamineum* and *U. bromivora* appeared to have several frameshifts introducing stop codons, suggesting that these sequences are pseudogenes. We reconstructed the phylogeny of the nucleotide sequences after removing ambiguously aligned regions (Figure 2). The resulting tree

shows that the closest relative of the *UMAG_11064* gene is the intronic sequence from *S. reilianum*.

As the GC profile of *UMAG_11064* suggests that the upstream region also has a mitochondrial origin (Figure 1b), we performed a codon alignment of the 5' region with the full intron sequences of *S. reilianum*, *T. indica*, and *T. walkeri* as well as the sequence of I-AbIII-P from *A. bisporus* in order to search for putative traces of the activity domain of the HE (Figure 3). We found that the intergenic region between *UMAG_11065* and *UMAG_11064* is similar to the activity domain of other GIY-YIG HE, and contains remnants of the former active site of the type GYY-YIG (Figure 3). Compared to I-AbIII-P and homologous sequences in *Tilletia*, however, a frameshift mutation has occurred in the active site (a 7 bp deletion). The predicted gene model for *UMAG_11064* starts at a conserved methionine position, 14 amino-acids downstream of the former active site (Figure 3) and contains the *helix-turn-helix* DNA-binding domain of the original HE.

A branch model of codon sequence evolution was fitted to the codon sequence alignment of *UMAG_11064* and its identified homologs, with the exception of the putative pseudogenes from *S. scitamineum* and *U. bromivora*. The proteins appear to be evolving under purifying selection ($dN/dS < 1$) on most branches of the tree, with the exception of the branches leading to the two *Tilletia* sequences, as well as the branch leading to *UMAG_11064* (Figure 2). We note, however, that the codeml program suffers from convergence issues on these particular branches, as witnessed by the variance in final estimates after 10 independent runs (Figure 2). Such convergence failures likely result from the branch model being an overparameterized model, here fitted on relatively short sequences. To further assess whether the high dN/dS ratio measured in the *UMAG_11064* gene could be explained both by relaxed purifying selection or positive selection, we fitted a branch-site model allowing specifically for

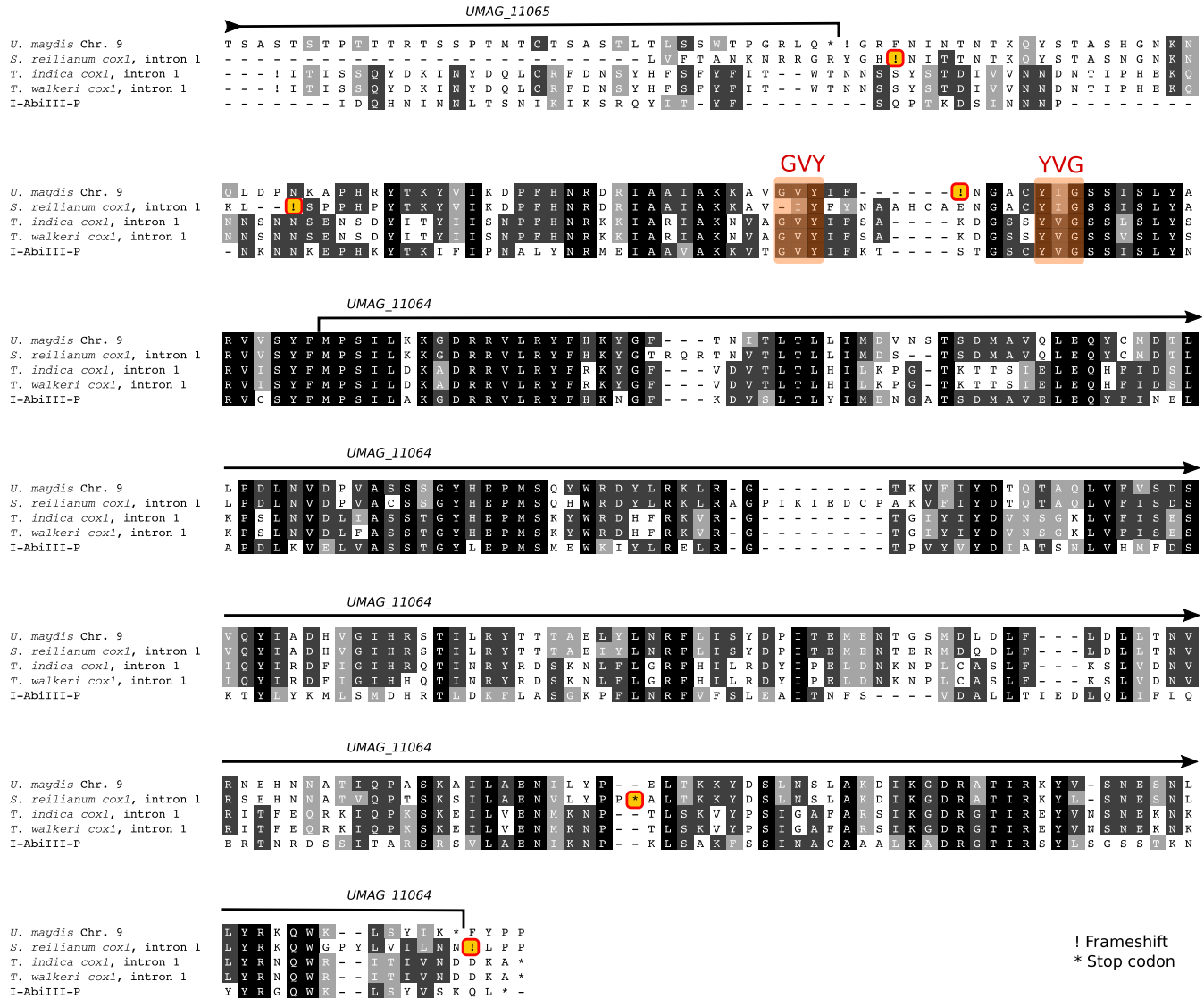


FIGURE 3 Alignment of UMAG_11064 and its upstream sequence with intron 1 from the *cox1* gene of *Sporisorium reilianum*, *Tilletia indica*, and *Tilletia walkeri*, as well as the coding sequence of the *Agaricus bisporus* HE. Shading indicates the level of amino-acid conservation, showing conserved residues (in black) and residues with similar biochemical properties (grayscale). Amino-acids noted as “X” have incomplete codons due to frameshifts. Highlighted exclamation marks denote inferred frameshifts and “*” characters stop codons. The location of the active site of the HE (GVVY-YVG) is highlighted

sites in the UMAG_11064 gene to evolve under positive selection (foreground branch), which we contrasted with a null model where all sites evolve under purifying selection or neutral evolution. The likelihood ratio test was not significant (p -value = .1585), even after removing the *Tilletia* sequences (p -values = .2183), and does not reject the hypothesis that the UMAG_11064 gene is evolving under a nearly neutral scenario. The higher dN/dS in the branch leading to UMAG_11064 might, therefore, be the result of relaxed purifying selection.

Altogether, these results suggest that UMAG_11064 is a former HE that inserted into the nuclear genome from the mitochondrion, was then inactivated by a deletion in its active site and acquired a new start codon, allowing it to code for a protein sequence with the former nucleotide-binding domain of the HE.

3.3 | The UMAG_11064 gene is similar to an intronic mitochondrial sequence of *S. reilianum*

The closest homologous sequence of UMAG_11064 was found in the first intron of the *cox1* gene of the smut fungus *S. reilianum* while this sequence was absent in the mitochondrial genome of *U. maydis*. The *cox1* genes of *S. reilianum* and *U. maydis* both have eight introns, of which only seven are homologous in position and sequence (Figure 4). *S. reilianum* has one extra intron in position 1, while *U. maydis* has one extra intron in position 6. In *U. maydis*, all introns but the sixth one are reported to contain a HEG. A blast search of this intron's sequence, however, revealed similarity with a homing endonuclease of type LAGLIDADG (Table S4). In *S. reilianum*, intron 1 (the putative precursor of UMAG_11064) and intron 2 are not annotated as containing

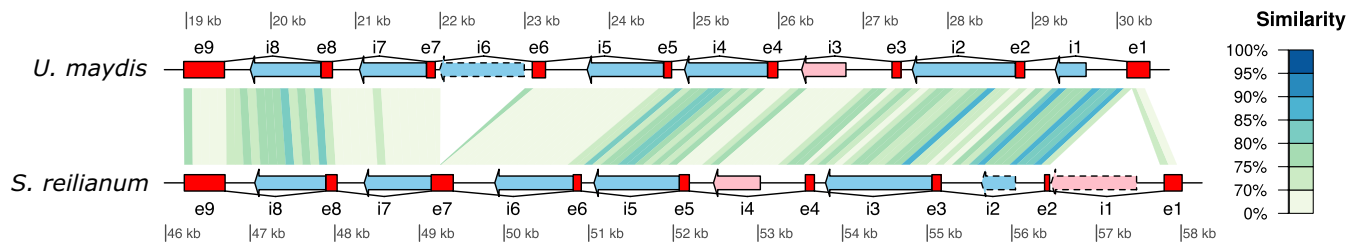


FIGURE 4 Intron structure of the *cox1* gene in *Ustilago maydis* and *Sporisorium reilianum*. Annotated HEs are indicated. Red boxes depict *cox1* exons, numbered from *e1* to *e9*. Introns are represented by connecting lines and numbered *i1* to *i8*. Arrows within introns show LAGLIDADG (light blue) and GIY-YIG HEs (pink). Dashed arrows correspond to HEGs inferred by blast search, while solid arrows correspond to the annotation from the GenBank files. Piecewise sequence similarity between *U. maydis* and *S. reilianum* is displayed with a color gradient

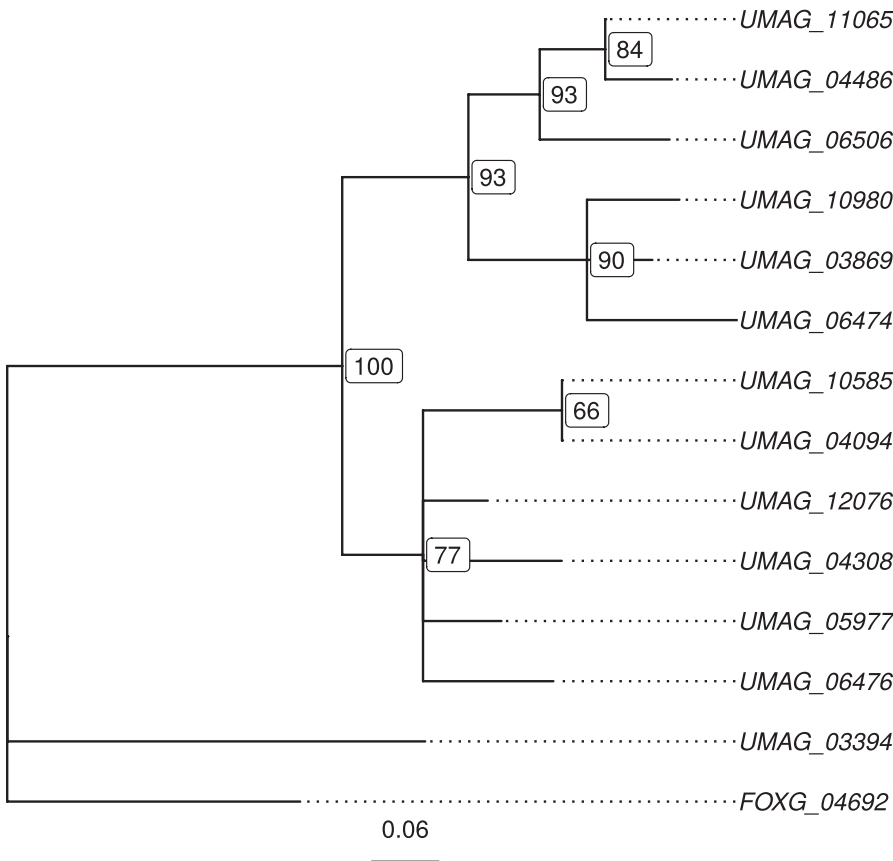


FIGURE 5 Maximum likelihood phylogeny of *UMAG_11065* *Ustilago maydis* paralogs together with the closest homolog from *F. oxysporum* (see Table 1). Node labels indicate support values as percentage. Nodes with support values lower than 60% have been collapsed

a HEG. Blast searches of the corresponding sequences, however, provided evidence for homology with a GIY-YIG HE (Table S5) and a LAGLIDADG HE, respectively (Table S6). Like many other species of fungi (Jalalzadeh et al., 2015; Pogoda et al., 2019; Stone et al., 2018), plants (Cho, Qiu, Kuhlman, & Palmer, 1998), and even animals (Fukami, Chen, Chiou, & Knowlton, 2007; Schuster et al., 2017), the *cox1* gene seems to be a hotspot of HEG-encoding introns in smut fungi.

Lastly, no homolog of intron 1 in *S. reilianum* was detected in the mitochondrial genome of *U. maydis*. A closer inspection showed that the ORF could be aligned with related HEs in other fungi (Figure 3). This alignment revealed an insertion of four amino-acids, a deletion of the first glycine residue in the active site plus several frameshifts at the beginning of the gene, which suggests that this gene has been altered and might not encode a functional HE any longer.

3.4 | *UMAG_11064* inserted into a gene encoding a RecQ helicase

In order to study the effect of the HEG insertion in the nuclear genome, we looked at the genomic environment of the *UMAG_11064* gene. Downstream of *UMAG_11064* is telomeric repeats, while the next upstream gene, *UMAG_11065*, is uncharacterized. A similarity search for *UMAG_11065* detected 13 paralogous sequences in the *U. maydis* genome (including one, *UMAG_12076*, on an unmapped contig), but only low-similarity matches in other sequenced smut fungi (see Methods). The closest nonsmut-related sequence comes from a gene from *F. oxysporum*. We inferred the evolutionary relationships between the 14 genes by reconstructing a maximum likelihood phylogenetic tree, and found that the *UMAG_11065* gene is closely related to *UMAG_04486*, located on chromosome 14

(Figure 5 and Table 1). The *UMAG_04486* gene, however, is predicted to be almost six times as long as *UMAG_11065*. We note that the downstream region of *UMAG_11064* does not show any similarity with the 3' part of the *UMAG_04486* gene, suggesting that the insertion of the HEG did not lead to the formation of an intron in the *UMAG_11065* gene, but rather to its truncation. A search for similar sequences of *UMAG_11065* and its relatives in public databases revealed homology with so-called RecQ helicases (Table S3), enzymes known to be involved in DNA repair and telomere expansion (Singh, Ghosh, Croteau, & Bohr, 2012). While this function is only predicted by homology, we note that all 12 chromosomal *recQ*-related genes are located very close to telomeres in *U. maydis* (Table 1), suggesting a role of these genes in telomere maintenance (Sánchez-Alonso & Guzmán, 1998). Lastly, we tested whether the truncation of *UMAG_11065* was followed by positively selected mutations in the remaining part of the gene. We inferred a dN/dS ratio equal to 0.342, which suggests that the *UMAG_11065* gene evolved mostly under purifying selection since divergence from the *UMAG_04486* gene. The insertion of *UMAG_11064*, therefore, was not followed by positive selection, or was too recent for sufficient positively selected substitutions to occur.

3.5 | *Ustilago maydis* populations show structural polymorphism in the telomeric region of chromosome 9

Because the *UMAG_11064* gene still displays a strong signature of its mitochondrial origin (codon usage and GC content), its transfer may have occurred recently. In order to provide a timeframe for the insertion event, we examined the structure of the genomic region of

the insertion in other *U. maydis* and *S. reilianum* isolates, as well as the structure of the *cox1* exons 1, 2, and 7. The regions that could be amplified and their corresponding sizes are listed in Figure S2, and the inferred genome organizations are summarized in Figure 6. The *UMAG_11064* gene is present in the FB1-derived strain SG200, as well as in the Holliday strains 518 and 521, but is absent in the nuclear and mitochondrial genome sequences of a recent *U. maydis* isolate from the US, strain 10-1, as well as from 5 Mexican isolates (I2, O2, P2, S5, and T6, Figure S2a). Conversely, the *UMAG_11072* gene, which is located further away from the telomere on the same chromosome arm, could be amplified in all strains. This positive control demonstrates that the lack of amplification of *UMAG_11064* in some strains is not due to any issue with the quality of the extracted DNA (Figure S2b). These results suggest that either the *UMAG_11064* gene was ancestral to all tested strains and subsequently lost in the Mexican and 10-1 strains, or it inserted in an ancestor of the three strains 518, 521 and SG200, after the divergence from other *U. maydis* strains, an event that occurred after the domestication of maize and the spread of the associated pathogen, 10,000 to 6,000 years ago (Munkacsí, Stoxen, & May, 2008). Moreover, all *U. maydis* strains possess intron 6 in the mitochondrial *cox1* gene, which is absent in *S. reilianum*. While the three *S. reilianum* strains tested carry intron 1, the most direct descendant of the progenitor of the HEs, it was absent in all *U. maydis* strains tested (Figures S2c–e and 6).

3.6 | Functional characterization

To shed light on the functional implication of the translocation of the HEG and subsequent mutations, we (a) assessed the expression profile of these genes and (b) generated a deletion strain and

TABLE 1 *UMAG_11065* paralogs in *Ustilago maydis*, together with a homolog from *Fusarium oxysporum* for comparison

| Gene | Chr/Scaffold/ Contig | Start | End | Length of Chr/ Scaffold/Contig | Number of introns | Length of protein | Relative position ^a (%) |
|------------|-------------------------|-----------|-----------|-----------------------------------|----------------------|----------------------|---------------------------------------|
| UMAG_06476 | Chromosome 3 | 1,641,500 | 1,642,057 | 1,642,070 | 0 | 185 | 99.98 |
| UMAG_06474 | Chromosome 3 | 1,639,598 | 1,640,203 | 1,642,070 | 0 | 201 | 99.87 |
| UMAG_06506 | Chromosome 7 | 951,043 | 954,234 | 957,188 | 5 | 983 | 99.52 |
| UMAG_10585 | Chromosome 4 | 883,585 | 884,046 | 884,984 | 0 | 153 | 99.87 |
| UMAG_11065 | Chromosome 9 | 1886 | 1,263 | 733,962 | 0 | 207 | 0.21 |
| UMAG_03394 | Chromosome 9 | 8,836 | 5,960 | 733,962 | 0 | 958 | 1.01 |
| UMAG_03869 | Chromosome 10 | 687,301 | 690,648 | 692,354 | 7 | 937 | 99.51 |
| UMAG_04094 | Chromosome 11 | 688,670 | 689,965 | 690,620 | 0 | 431 | 99.81 |
| UMAG_04486 | Chromosome 14 | 605,233 | 609,089 | 611,467 | 2 | 1,175 | 99.30 |
| UMAG_04308 | Chromosome 14 | 1,241 | 87 | 611,467 | 0 | 384 | 0.11 |
| UMAG_05977 | Chromosome 20 | 523,510 | 523,884 | 523,884 | 0 | 124 | 99.96 ^b |
| UMAG_10980 | Chromosome 22 | 398,220 | 400,499 | 403,590 | 0 | 759 | 98.95 |
| UMAG_12076 | Contig 1.265 | 4,214 | 5,343 | 5,343 | 0 | 376 | 89.43 |
| FOXG_04692 | Supercontig 2.5 | 9,736 | 6,398 | 2,688,632 | 0 | 1,112 | 0.30 |

^aPosition reported to the length of the chromosome or contig.

^bN-terminal fragment only.

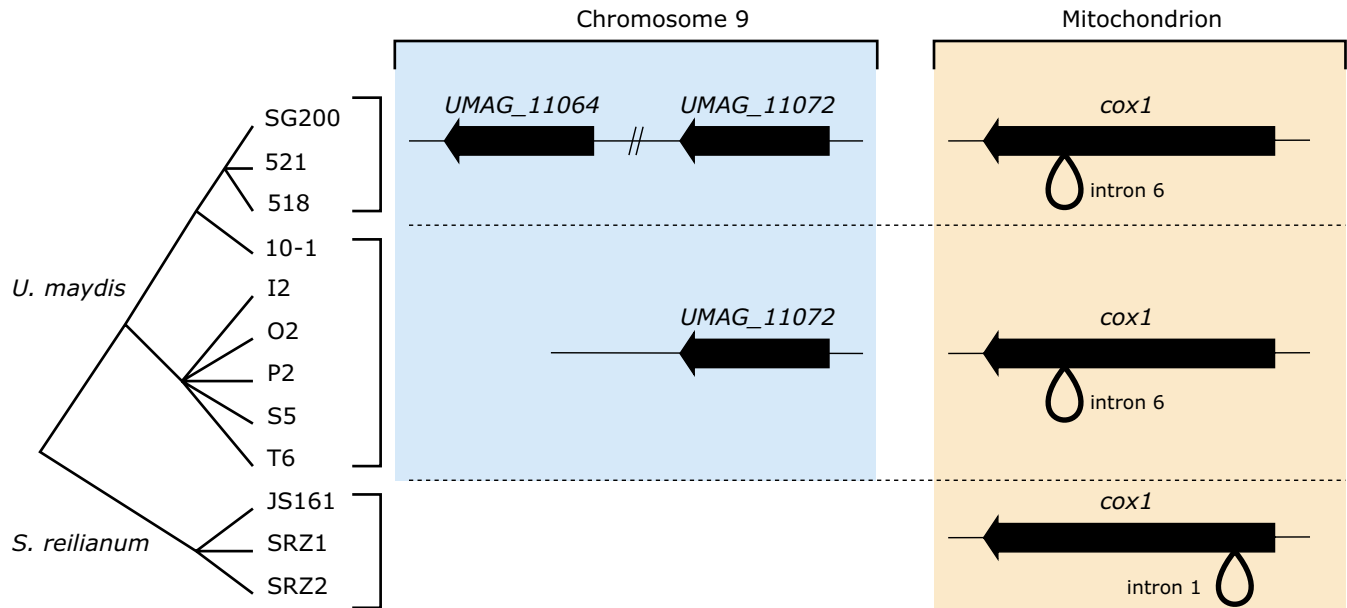


FIGURE 6 Presence of the *UMAG_11064* gene and structure of the *cox1* gene in several *Ustilago maydis* and *Sporisorium reilianum* strains, as assessed by PCR, together with their phylogeny. The *UMAG_11072* gene, located 90 kb downstream the *UMAG_11064* gene on chromosome 9, was used as a positive control. Strains 521 and 518 are two strains resulting from the same spore from a field isolate from the United States. SG200 is a genetically engineered strain derived from a cross between the 518 and 521 strains. Strain 10-1 is another field isolate from the United States. Strains I2, O2, P2, S5, and T6 from field isolates from Mexico

phenotyped it. For the expression analysis, we relied on a previously published RNASeq data set (Lanver et al., 2018), from which we extracted the expression profiles of genes in the telomeric region of chromosome 9 (Figure 7a). While the expression of *UMAG_11064* remained close to zero in the three replicates, expression of *UMAG_11065* increased during plant infection. The telomeric region was highly heterogeneous in terms of expression profile: While *UMAG_11066* and *UMAG_03393* did not show any significant level of expression, *UMAG_03392* was down-regulated at 12 hr postinfection, while *UMAG_03394*, another RecQ-encoding gene homologous to *UMAG_11065*, displayed constitutively high levels of expression (Figure 7a). All homologs of *UMAG_11065* show a significantly higher expression during infection (Tukey's post hoc test, false discovery rate of 5%, Figure 7b). The comparison of expression profiles revealed two main classes of genes (Figure 7c): highly expressed genes (upper group) and moderately expressed genes (lower group), to which *UMAG_11065* belongs. We further note that the differences in expression profiles do not mirror the protein sequence similarity of the genes (Mantel permutation test, p -value = .566).

To assess the functional role of *UMAG_11064* and *UMAG_11065*, these genes were simultaneously deleted in SG200, a solopathogenic haploid strain that can cause disease without a mating partner (Kämper et al., 2006) using a single-step gene replacement method (Kämper, 2004). Gene deletion was verified by Southern analysis (Figure S3). Virulence assays conducted in triplicate revealed no statistically different symptoms of the double deletion strain, SG200 Δ 11065 Δ 11064, compared to SG200 in infected maize plants (Figure 8a, chi-squared test, p -value = .453). Since RecQ helicases contribute to dealing with replication stress (Kojic & Holloman, 2012),

we also determined the sensitivity of the mutant to various stressors including UV, hydroxyurea, and Congo Red. (Figure 8b). We report that the deletion strain shows increased sensitivity to cell wall stress induced by Congo Red and increased resistance to UV stress. Since *UMAG_11064* does not show any detectable level of expression, we hypothesize that the deletion of *UMAG_11065* is responsible for this phenotype.

4 | DISCUSSION

The codon usage and GC content of the *UMAG_11064* gene, as well as its similarity to known mitochondrial HEGs, point at a recent transfer into the nuclear genome of *U. maydis*. Moreover, the precursor of this gene is absent from the mitochondrial genome of this species. Two possible scenarios can explain this pattern, which we detail below.

The first scenario involves a transfer of the gene to the nuclear genome followed by a loss of the mitochondrial copy (Figure 9). Under this scenario, the mitochondrial HEG was present in the *U. maydis* ancestor. Two evolutionary events are invoked: the insertion of the HEG into the nuclear genome, on the one hand, creating a HEG⁺ genotype at the nuclear locus (designated [HEG⁺]_{nuc}), and the loss of the mitochondrial copy, creating a HEG⁻ genotype at the mitochondrial locus (designated [HEG⁻]_{mit}). These two events may have happened at distinct time points, but, under this scenario, the former cannot have happened after the fixation of the [HEG⁻]_{mit} genotype in the population. The [HEG⁺]_{nuc}/ [HEG⁻]_{mit} genotype could be generated by a cross between two individuals,

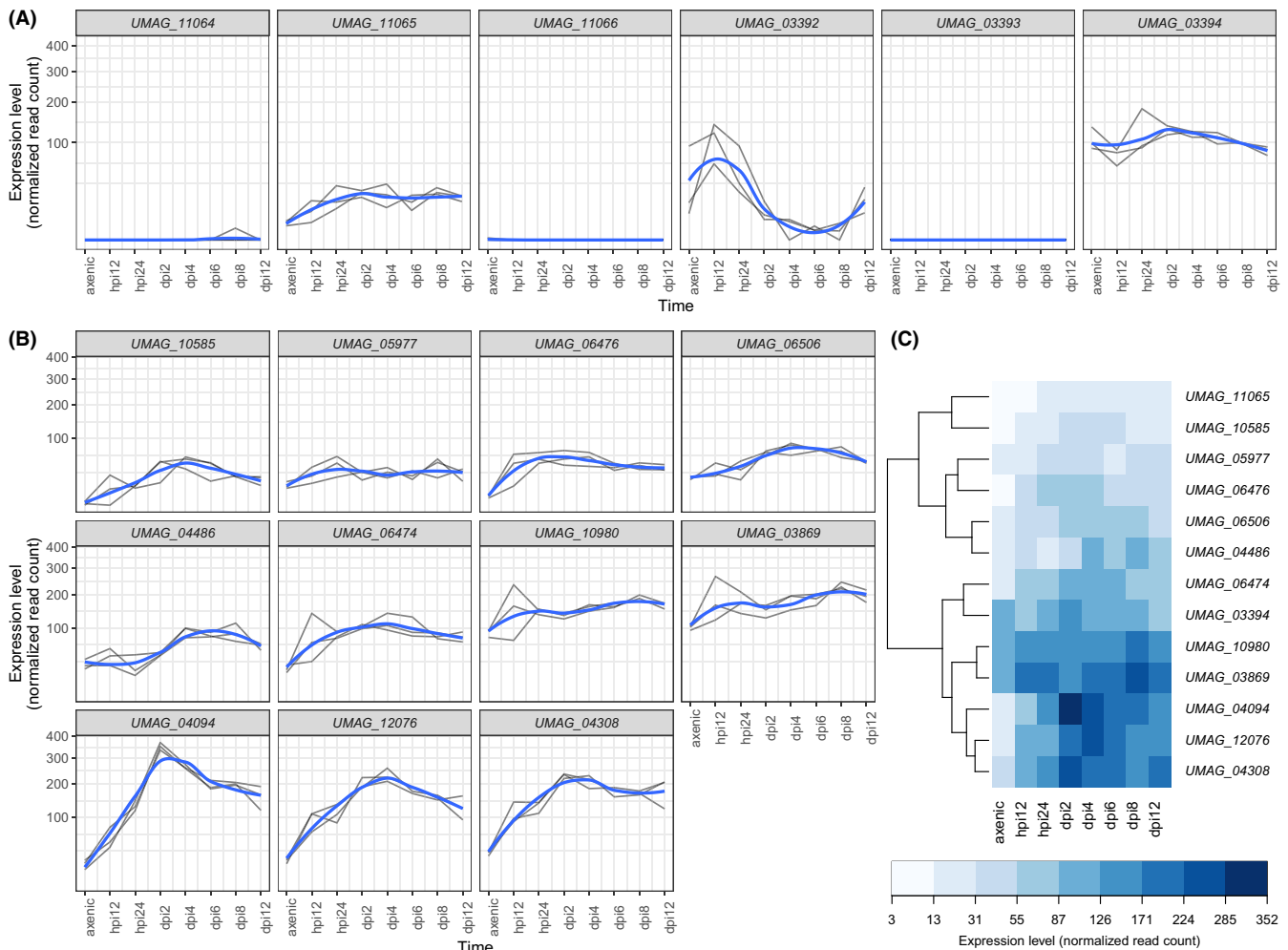


FIGURE 7 Patterns of gene expression for *UMAG_11064* and *UMAG_11065*, together with neighboring and homologous genes. (a) Gene expression profiles for genes in the chromosome 9 telomeric region (as depicted on Figure 1b). Straight lines represent three independent replicates, while the blue curve depicts the smoothed conditional mean computed using the LOESS method. (b) Gene expression profiles for the *UMAG_11065* homologs (Figure 5). Legends as in (a). (c) Clustering of the *UMAG_11065* homologs based on their averaged expression profile (see Methods). Hpi, hours postinfection; Dpi, days postinfection

one $[HEG^+]_{nuc}$ and the other $[HEG^-]_{mit}$, given that mitochondria are uniparentally inherited in *U. maydis* (Basse, 2010). Importantly, the segregation of the $[HEG^+]_{nuc}$ and $[HEG^-]_{mit}$ variants could be purely neutral and driven by genetic drift only. In case the $[HEG^+]_{nuc}$ allele contained a recognition sequence of the HE, the $[HEG^+]_{nuc}$ allele may have initially benefited from a genetic drive effect. Any putative selective advantage/disadvantage of the $[HEG^+]_{nuc}$ or $[HEG^-]_{mit}$ alleles may have favored their fixation, or on the contrary, acted against it.

In the second scenario, the mitochondrial HEG was not ancestral to *U. maydis*, but was horizontally transferred from *S. reilianum* (or a related species). The high similarity of the *UMAG_11064* gene to the *S. reilianum* mitochondrial HEG (Figure 3) supports this hypothesis, given the relatively high nucleotide divergence between the two species, which diverged around 20 My ago (Schweizer et al., 2018). We note, however, that intronic HEGs have also been reported to show reduced nucleotide substitution rates, which can potentially explain their comparatively low divergence (Jalalzadeh

et al., 2015). Group I introns have been reported to be highly mobile and to undergo frequent horizontal gene transfers (HGT) within metazoans (Schuster et al., 2017), plants (Sanchez-Puerta, Cho, Mower, Alverson, & Palmer, 2008), and fungi (Férandon et al., 2010; Jalalzadeh et al., 2015). Group I introns in metazoans and plants are also thought to originate from a fungal donor (Sanchez-Puerta et al., 2008, 2011; Schuster et al., 2017). In this respect, a transfer from another *cox1* intron-carrying smut fungus to *U. maydis* is not unlikely, given that *U. maydis* and *S. reilianum* share the same host, and that hybridization between smut species has been reported (Boidin, 1986; Fischer, 1957).

The insertion of the HEG in the nuclear genome poses the question of the underlying mechanisms, independently of the origin of the HEG. First, the HEG could be encoding a fully functional HE and the $[HEG^-]_{nuc}$ allele contained a HE recognition sequence. Under this scenario, the insertion event was a homing event and the inactivation of the HEG occurred after the insertion. Therefore, several generations must have passed since the insertion event in order for the

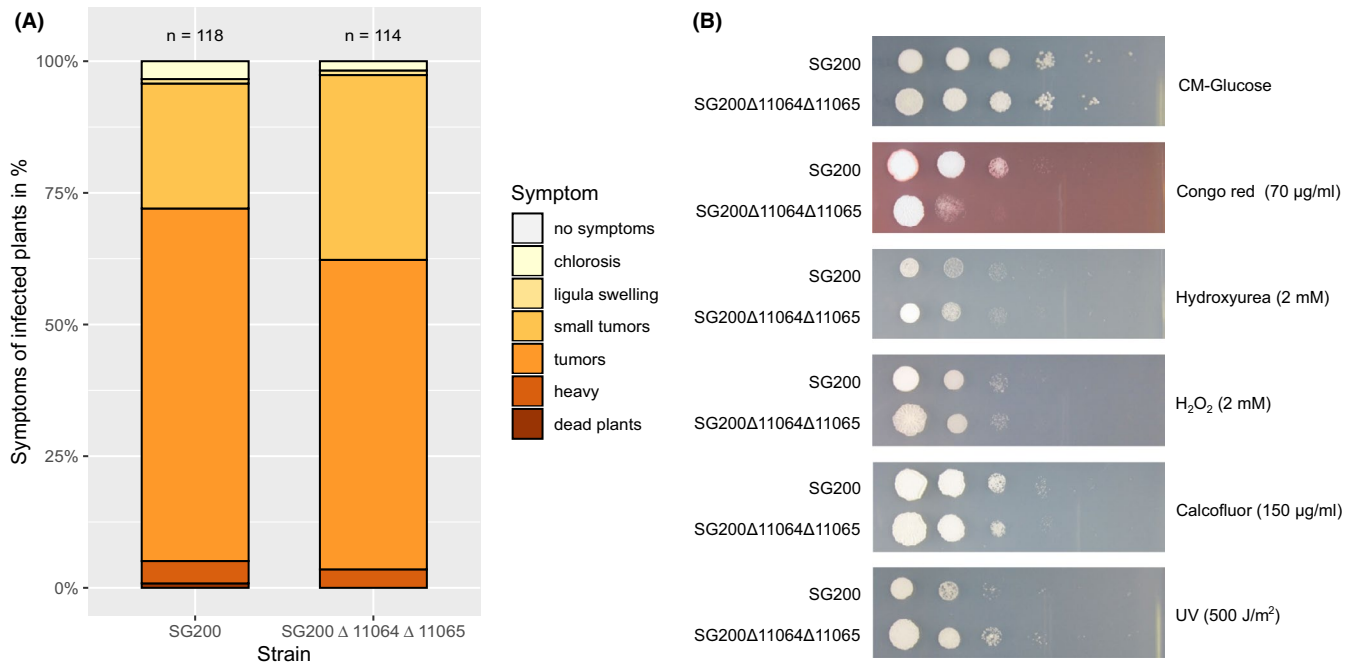


FIGURE 8 Phenotype assessment of the double deletion strain. (a) The simultaneous deletion of *UMAG_11064* and *UMAG_11065* does not affect virulence. Maize seedlings were infected with the indicated strains. Disease symptoms were scored at 12 dpi according to Kämper et al. (2006) using the color code depicted on the right. Colors reflect the degree of severity, from brown-red (severe) to light yellow (mild). Data represent mean of $n = 3$ biologically independent experiments. Total numbers of infected plants are indicated above the respective columns. (b) Stress assay of the double deletion strain (Δ 11064 Δ 11065), lacking both genes *UMAG_11064* and *UMAG_11065*, compared to the parental SG200 strain. Assays were repeated at least three times with comparable results

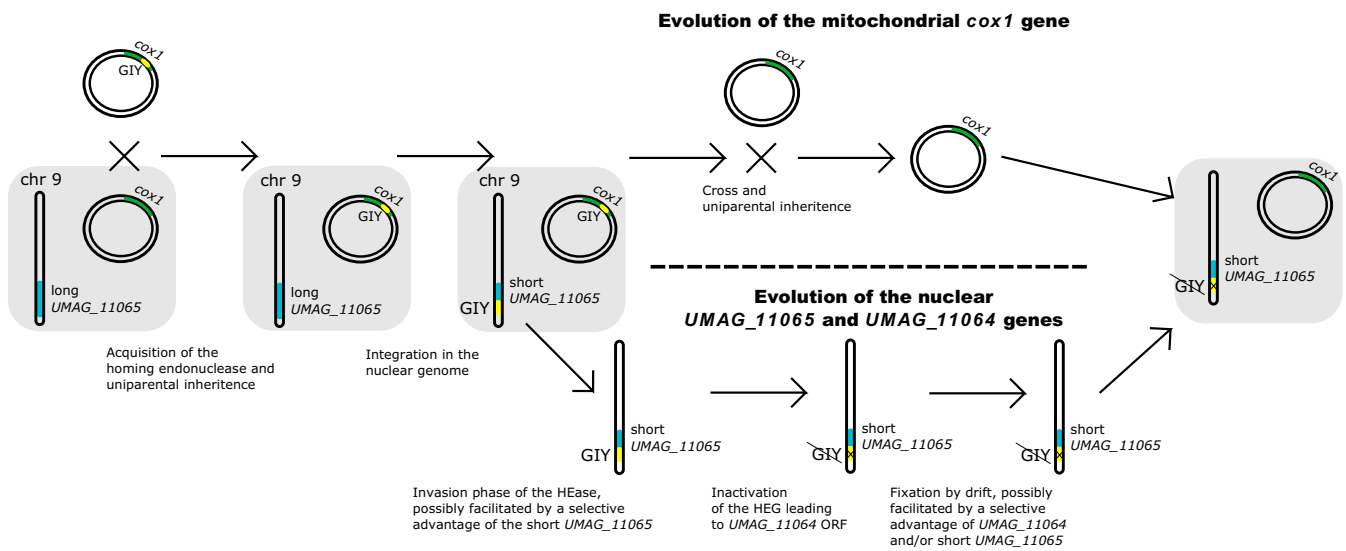


FIGURE 9 Possible evolutionary scenario recapitulating the events leading to the formation of the *UMAG_11064* and *UMAG_11065* *Ustilago maydis* genes. Importantly, each step in this model may have occurred by genetic drift alone. Positive selection may have favored—but is not required to explain—the spread of the nuclear and/or the loss of the mitochondrial HEGs

inactivating mutations to occur. Alternatively, the [HEG]_{nuc} allele might not have contained a recognition sequence and the insertion of the HEG occurred by an unknown mechanism. Lastly, a possibility is that the inserted sequence already encoded an inactivated HE, which was then inserted by an unknown mechanism. In this latter case, the insertion could have occurred very recently, possibly a few generations in the past.

Homing endonuclease genes are found in eukaryotic nuclei but are usually restricted to small and large ribosomal RNA subunit genes (Dunin-Horkawicz, Feder, & Bujnicki, 2006; Lambowitz & Belfort, 1993). While transfer of DNA segments and functional genes from organellar genomes to the nucleus is well documented (Fuentes, Karcher, & Bock, 2012; Lloyd & Timmis, 2011; Sun & Callis, 1993; Thorsness & Weber, 1996), established examples of

HEG insertions at other genomic locations than rRNA genes are very scarce. Several questions remain unanswered regarding the mechanisms of insertion into the nuclear genome, providing it happened as a homing event. For the event to happen, a template sequence containing the HEG is required for repairing the break initiated by the HE. This template must have, therefore, "leaked" from the mitochondrion. The recognition motif of the original *UMAG_11064* HEG is unknown, and the very short flanking regions surrounding the insertion site do not allow any comparison with known motifs, preventing further conclusions to be made regarding the nature of the insertion event of *UMAG_11064*.

Interestingly, Louis and Haber (Louis & Haber, 1991) reported a similar transfer of a HEG into a telomeric region of *Saccharomyces cerevisiae*. The authors argue that signatures of such insertion could be found because (a) it had no deleterious effect and (b) the occurrence of heterologous recombination between telomeres favors the maintenance of elements that would otherwise be lost. Contrasting with this result, the insertion of the GIY-YIG HEG that inserted into the ancestor of the *UMAG_11065* gene potentially had non-neutral effects, resulting in an expressed truncated protein. Several mutations were found within the active site of the inserted HEG that led to the *UMAG_11064* gene, suggesting that the encoded protein is unlikely to act as a HE any longer. However, a putative alternative start codon was detected, downstream the active site, followed by an uninterrupted peptide sequence containing the helix-turn-helix binding domain of the original HE. Furthermore, we could not detect any significant level of expression of the *UMAG_11064* gene in various laboratory conditions. Comparative sequence analysis further suggests that *UMAG_11064* is evolving under relaxed purifying selection, indicating that it might be undergoing pseudogenization. These results, therefore, suggest that the *UMAG_11064* gene is not functional. However, as this mitochondrial HEG inserted into a nuclear *U. maydis* gene, it might have had phenotypic consequences not directly due to the HEG gene itself. The *UMAG_11065* gene appeared truncated by the HEG insertion, which removed the C-terminal part of the encoded protein, a likely RecQ helicase, and the truncated *UMAG_11065* is expressed during infection. The truncation likely did not have a strong negative impact, possibly because of the existence of multiple potentially functionally redundant paralogs of *UMAG_11065*, including on the same telomeric region of chromosome 9, with *UMAG_03394* being located 4 genes upstream (Table 1). While we were unable to detect a contribution to virulence, our results point at a putative role of the truncated RecQ helicase into stress tolerance, as its deletion increases resistance to UV radiation but makes the fungus more susceptible to cell wall stress, at least under laboratory conditions. How the truncated *UMAG_11065* RecQ helicase could improve coping with cell wall stress and increases the sensitivity to UV simultaneously, however, remains to be investigated, as well as the potential fitness benefit or cost of these phenotypes. Furthermore, a possibility remains that the observed phenotype of *UMAG_11065* is ancestral and not due to the truncation itself, which could be neutral. In order to elucidate the putative adaptive role of the truncation of *UMAG_11065*, knowledge

of the ancestral, nontruncated *UMAG_11065* allele is needed, as well as its distribution in natural populations.

5 | CONCLUSIONS

In this study, we report instances of two stages of the life cycle of HEGs. Intron 1 of the mitochondrial *cox1* gene of *S. reilianum* was shown to contain a degenerated GIY-YIG HEG, while the homologous position in the *U. maydis* gene displays no intron. Besides, in the telomeric region of chromosome 9 of the nuclear genome of *U. maydis*, we found evidence of a recent insertion of a very similar GIY-YIG HEG. Phenotypic assay of the mutant strain containing a double deletion of the HEG and the helicase gene where it inserted reveals enhanced stress sensitivity in vitro. The absence of a GIY-YIG HEG in any field isolates of *U. maydis* sequenced so far, however, suggests that either the mutation was lost in natural populations and only maintained under laboratory conditions, or that it is only present in a so far unsampled population. The *UMAG_11064* gene offers a snapshot of evolution taken soon after a mutation event occurred. As such, it can provide insights into the mechanisms of HEG mobility and horizontal transfer. These results further demonstrate that HEGs can generate genetic diversity not only via their duplication, but also by drastically modifying the local genome architecture where they insert.

ACKNOWLEDGMENTS

We thank all members of our groups for stimulating discussions. We are grateful to Georgiana May and Octavio Paredes-López for providing field isolates of *U. maydis* from the United States and Mexico, respectively. We acknowledge the generous support by the Max Planck Society. This manuscript has been peer-reviewed and recommended by Peer Community in Evolutionary Biology (<https://doi.org/10.24072/pci.evolbiol.100101>). The present version of this manuscript greatly benefited from the comments of the recommender, Sylvain Charlat, and the two reviewers, Jan Engelstaedter and Yannick Wurm.

CONFLICT OF INTEREST

None declared.

AUTHOR CONTRIBUTION

Julien Y. Dutheil: Conceptualization (equal); Data curation (lead); Formal analysis (lead); Investigation (lead); Methodology (equal); Visualization (lead); Writing-original draft (lead); Writing-review & editing (equal). **Karin Münch:** Data curation (supporting); Formal analysis (supporting); Investigation (supporting); Methodology (supporting); Writing-review & editing (supporting). **Klaas Schotanus:** Methodology (supporting); Writing-review & editing (supporting). **Eva H. Stukenbrock:** Investigation (supporting); Methodology (supporting); Writing-review & editing (equal). **Regine Kahmann:** Conceptualization (equal); Funding acquisition (lead); Methodology (equal); Resources (lead); Writing-review & editing (equal).

DATA AVAILABILITY STATEMENT

Data sets and scripts used to conduct the phylogenetic and statistical analyses, as well as R code used to generate figures 1, 3, 4, 5, and 6 are available at https://gitlab.gwdg.de/molsysevol/umag_11064 and Zenodo (<https://doi.org/10.5281/zenodo.3984974>).

ORCID

Julien Y. Dutheil  <https://orcid.org/0000-0001-7753-4121>

Karin Münch  <https://orcid.org/0000-0002-7437-6823>

Klaas Schotanus  <https://orcid.org/0000-0002-0974-2882>

Eva H. Stukenbrock  <https://orcid.org/0000-0001-8590-3345>

Regine Kahmann  <https://orcid.org/0000-0001-7779-7837>

REFERENCES

- Altschul, S. F., Gish, W., Miller, W., Myers, E. W., & Lipman, D. J. (1990). Basic local alignment search tool. *Journal of Molecular Biology*, 215, 403–410. [https://doi.org/10.1016/S0022-2836\(05\)80360-2](https://doi.org/10.1016/S0022-2836(05)80360-2)
- Anisimova, M., & Gascuel, O. (2006). Approximate likelihood-ratio test for branches: A fast, accurate, and powerful alternative. *Systematic Biology*, 55, 539–552. <https://doi.org/10.1080/10635150600755453>
- Ast, J., Stiebler, A. C., Freitag, J., & Böcker, M. (2013). Dual targeting of peroxisomal proteins. *Frontiers in Physiology*, 4, 297. <https://doi.org/10.3389/fphys.2013.00297>
- Banuett, F., & Herskowitz, I. (1989). Different alleles of *Ustilago maydis* are necessary for maintenance of filamentous growth but not for meiosis. *Proceedings of the National Academy of Sciences of the United States of America*, 86, 5878–5882. <https://doi.org/10.1073/pnas.86.15.5878>
- Barzel, A., Obolski, U., Gogarten, J. P., Kupiec, M., & Hadany, L. (2011). Home and away—the evolutionary dynamics of homing endonucleases. *BMC Evolutionary Biology*, 11, 324. <https://doi.org/10.1186/1471-2148-11-324>
- Basse, C. W. (2010). Mitochondrial inheritance in fungi. *Current Opinion in Microbiology*, 13, 712–719. <https://doi.org/10.1016/j.mib.2010.09.003>
- Belfort, M., & Roberts, R. J. (1997). Homing endonucleases: Keeping the house in order. *Nucleic Acids Research*, 25, 3379–3388. <https://doi.org/10.1093/nar/25.17.3379>
- Boidin, J. (1986). Intercompatibility and the species concept in the saprobic basidiomycotina. *Mycotaxon*, XXVI, 319–336.
- Charif, D., Thioulouse, J., Lobry, J. R., & Perrière, G. (2005). Online synonymous codon usage analyses with the ade4 and seqinR packages. *Bioinformatics*, 21, 545–547. <https://doi.org/10.1093/bioinformatics/bti037>
- Chevalier, B. S., & Stoddard, B. L. (2001). Homing endonucleases: Structural and functional insight into the catalysts of intron/intron mobility. *Nucleic Acids Research*, 29, 3757–3774. <https://doi.org/10.1093/nar/29.18.3757>
- Cho, Y., Qiu, Y.-L., Kuhlman, P., & Palmer, J. D. (1998). Explosive invasion of plant mitochondria by a group I intron. *Proceedings of the National Academy of Sciences of the United States of America*, 95, 14244–14249. <https://doi.org/10.1073/pnas.95.24.14244>
- Christianson, T. W., Sikorski, R. S., Dante, M., Shero, J. H., & Hieter, P. (1992). Multifunctional yeast high-copy-number shuttle vectors. *Gene*, 110, 119–122. [https://doi.org/10.1016/0378-1119\(92\)90454-W](https://doi.org/10.1016/0378-1119(92)90454-W)
- Cock, P. J. A., Antao, T., Chang, J. T., Chapman, B. A., Cox, C. J., Dalke, A., ... de Hoon, M. J. L. (2009). Biopython: Freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics*, 25, 1422–1423. <https://doi.org/10.1093/bioinformatics/btp163>
- Derbyshire, V., Kowalski, J. C., Dansereau, J. T., Hauer, C. R., & Belfort, M. (1997). Two-domain structure of the td intron-encoded endonuclease I-TevI correlates with the two-domain configuration of the homing site. *Journal of Molecular Biology*, 265, 494–506. <https://doi.org/10.1006/jmbi.1996.0754>
- Djamei, A., & Kahmann, R. (2012). *Ustilago maydis*: Dissecting the molecular interface between pathogen and plant. *PLoS Path*, 8, e1002955. <https://doi.org/10.1371/journal.ppat.1002955>
- Dong, S., Raffaele, S., & Kamoun, S. (2015). The two-speed genomes of filamentous pathogens: Waltz with plants. *Current Opinion in Genetics & Development*, 35, 57–65. <https://doi.org/10.1016/j.gde.2015.09.001>
- Dujon, B., Belfort, M., Butow, R. A., Jacq, C., Lemieux, C., Perlman, P. S., & Vogt, V. M. (1989). Mobile introns: Definition of terms and recommended nomenclature. *Gene*, 82, 115–118. [https://doi.org/10.1016/0378-1119\(89\)90035-8](https://doi.org/10.1016/0378-1119(89)90035-8)
- Dunin-Horkawicz, S., Feder, M., & Bujnicki, J. M. (2006). Phylogenomic analysis of the GYI-YIG nuclease superfamily. *BMC Genomics*, 7, 98. <https://doi.org/10.1186/1471-2164-7-98>
- Dutheil, J. Y., Mannhaupt, G., Schweizer, G., Sieber, C. M. K., Münsterkötter, M., Güldener, U., ... Kahmann, R. (2016). A tale of genome compartmentalization: The evolution of virulence clusters in smut fungi. *Genome Biology and Evolution*, 8, 681–704. <https://doi.org/10.1093/gbe/evw026>
- Eyre-Walker, A., & Keightley, P. D. (2007). The distribution of fitness effects of new mutations. *Nature Reviews Genetics*, 8, 610–618. <https://doi.org/10.1038/nrg2146>
- Férandon, C., Moukha, S., Callac, P., Benedetto, J.-P., Castroviejo, M., & Barroso, G. (2010). The *Agaricus bisporus* cox1 Gene: The longest mitochondrial gene and the largest reservoir of mitochondrial Group I introns. *PLoS One*, 5. <https://doi.org/10.1371/journal.pone.0014048>
- Fischer, C. (1957). *Biology and control of smut fungi*. Toronto, ON: John Wiley & Sons Canada, Limited.
- Fuentes, I., Karcher, D., & Bock, R. (2012). Experimental reconstruction of the functional transfer of intron-containing plastid genes to the nucleus. *Current Biology*, 22, 763–771. <https://doi.org/10.1016/j.cub.2012.03.005>
- Fukami, H., Chen, C. A., Chiou, C.-Y., & Knowlton, N. (2007). Novel group I introns encoding a putative homing endonuclease in the mitochondrial cox1 gene of Scleractinian corals. *Journal of Molecular Evolution*, 64, 591–600. <https://doi.org/10.1007/s00239-006-0279-4>
- Gladyshev, E. (2017). Repeat-induced point mutation and other genome defense mechanisms in fungi. *Microbiology Spectrum*, 5, 687–699. <https://doi.org/10.1128/microbiolspec.FUNK-0042-2017>
- Goddard, M. R., & Burt, A. (1999). Recurrent invasion and extinction of a selfish gene. *Proceedings of the National Academy of Sciences of the United States of America*, 96, 13880–13885. <https://doi.org/10.1073/pnas.96.24.13880>
- Gogarten, J. P., & Hilario, E. (2006). Inteins, introns, and homing endonucleases: Recent revelations about the life cycle of parasitic genetic elements. *BMC Evolutionary Biology*, 6, 94.
- Gouy, M., Guindon, S., & Gascuel, O. (2010). SeaView version 4: A multiplatform graphical user interface for sequence alignment and phylogenetic tree building. *Molecular Biology and Evolution*, 27, 221–224. <https://doi.org/10.1093/molbev/msp259>
- Grandaubert, J., Lowe, R. G. T., Soyer, J. L., Schoch, C. L., Van de Wouw, A. P., Fudal, I., ... Rouxel, T. (2014). Transposable element-assisted evolution and adaptation to host plant within the *Leptosphaeria maculans*-*Leptosphaeria biglobosa* species complex of fungal pathogens. *BMC Genomics*, 15, 891. <https://doi.org/10.1186/1471-2164-15-891>
- Guindon, S., Dufayard, J.-F., Lefort, V., Anisimova, M., Hordijk, W., & Gascuel, O. (2010). New algorithms and methods to estimate maximum-likelihood phylogenies: Assessing the performance of PhyML 3.0. *Systematic Biology*, 59, 307–321. <https://doi.org/10.1093/sysbio/syq010>
- Guy, L., Kultima, J. R., & Andersson, S. G. E. (2010). genoPlotR: Comparative gene and genome visualization in R. *Bioinformatics*, 26, 2334–2335. <https://doi.org/10.1093/bioinformatics/btq413>

- Jalalzadeh, B., Saré, I. C., Férandon, C., Callac, P., Farsi, M., Savoie, J.-M., & Barroso, G. (2015). The intraspecific variability of mitochondrial genes of *Agaricus bisporus* reveals an extensive group I intron mobility combined with low nucleotide substitution rates. *Current Genetics*, *61*, 87–102. <https://doi.org/10.1007/s00294-014-0448-8>
- Kaessmann, H. (2010). Origins, evolution, and phenotypic impact of new genes. *Genome Research*, *20*, 1313–1326. <https://doi.org/10.1101/gr.101386.109>
- Kämper, J. (2004). A PCR-based system for highly efficient generation of gene replacement mutants in *Ustilago maydis*. *Molecular Genetics and Genomics*, *271*, 103–110. <https://doi.org/10.1007/s00438-003-0962-8>
- Kämper, J., Kahmann, R., Bölker, M., Ma, L.-J., Brefort, T., Saville, B. J., ... Birren, B. W. (2006). Insights from the genome of the biotrophic fungal plant pathogen *Ustilago maydis*. *Nature*, *444*, 97–101. <https://doi.org/10.1038/nature05248>
- Katoh, K., Misawa, K., Kuma, K., & Miyata, T. (2002). MAFFT: A novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Research*, *30*, 3059–3066. <https://doi.org/10.1093/nar/gkf436>
- Khrunyk, Y., Münch, K., Schipper, K., Lupas, A. N., & Kahmann, R. (2010). The use of FLP-mediated recombination for the functional analysis of an effector gene family in the biotrophic smut fungus *Ustilago maydis*. *New Phytologist*, *187*, 957–968. <https://doi.org/10.1111/j.1469-8137.2010.03413.x>
- Kojic, M., & Holloman, W. K. (2012). Brh2 domain function distinguished by differential cellular responses to DNA damage and replication stress. *Molecular Microbiology*, *83*, 351–361. <https://doi.org/10.1111/j.1365-2958.2011.07935.x>
- Kondrashov, F. A., & Kondrashov, A. S. (2010). Measurements of spontaneous rates of mutations in the recent past and the near future. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, *365*, 1169–1176. <https://doi.org/10.1098/rstb.2009.0286>
- Krombach, S., Reissmann, S., Kreibich, S., Bochen, F., & Kahmann, R. (2018). Virulence function of the *Ustilago maydis* sterol carrier protein 2. *New Phytologist*, *220*, 553–566.
- Lambowitz, A. M., & Belfort, M. (1993). Introns as mobile genetic elements. *Annual Review of Biochemistry*, *62*, 587–622. <https://doi.org/10.1146/annurev.bi.62.070193.003103>
- Lanver, D., Müller, A. N., Happel, P., Schweizer, G., Haas, F. B., Franitz, M., ... Kahmann, R. (2018). The biotrophic development of *Ustilago maydis* studied by RNA-Seq analysis. *The Plant Cell*, *30*, 300–323.
- Laurie, J. D., Ali, S., Linning, R., Mannhaupt, G., Wong, P., Güldener, U., ... Schirawski, J. (2012). Genome comparison of barley and maize smut fungi reveals targeted loss of RNA silencing components and species-specific presence of transposable elements. *The Plant Cell*, *24*, 1733–1745. <https://doi.org/10.1105/tpc.112.097261>
- Le, S. Q., & Gascuel, O. (2008). An improved general amino acid replacement matrix. *Molecular Biology and Evolution*, *25*, 1307–1320. <https://doi.org/10.1093/molbev/msn067>
- Lloyd, A. H., & Timmis, J. N. (2011). The origin and characterization of new nuclear genes originating from a cytoplasmic organellar genome. *Molecular Biology and Evolution*, *28*, 2019–2028. <https://doi.org/10.1093/molbev/msr021>
- Louis, E. J., & Haber, J. E. (1991). Evolutionarily recent transfer of a group I mitochondrial intron to telomere regions in *Saccharomyces cerevisiae*. *Current Genetics*, *20*, 411–415. <https://doi.org/10.1007/BF00317070>
- Lynch, M. (2007). *The origins of genome architecture*. Sunderland, MA: Sinauer Associates.
- Lynch, M., Sung, W., Morris, K., Coffey, N., Landry, C. R., Dopman, E. B., ... Thomas, W. K. (2008). A genome-wide view of the spectrum of spontaneous mutations in yeast. *Proceedings of the National Academy of Sciences*, *105*, 9272–9277. <https://doi.org/10.1073/pnas.0803466105>
- Mewes, H. W., Ruepp, A., Theis, F., Rattei, T., Walter, M., Frishman, D., ... Antonov, A. (2011). MIPS: Curated databases and comprehensive secondary data resources in 2010. *Nucleic Acids Research*, *39*, D220–224. <https://doi.org/10.1093/nar/gkq1157>
- Möller, M., & Stukenbrock, E. H. (2017). Evolution and genome architecture in fungal plant pathogens. *Nature Reviews Microbiology*, *15*, 756–771. <https://doi.org/10.1038/nrmicro.2017.76>
- Munkacsy, A. B., Stoxen, S., & May, G. (2008). *Ustilago maydis* populations tracked maize through domestication and cultivation in the Americas. *Proceedings of the Royal Society B: Biological Sciences*, *275*, 1037–1046.
- Ohta, T. (2000). Evolution of gene families. *Gene*, *259*, 45–52. [https://doi.org/10.1016/S0378-1119\(00\)00428-5](https://doi.org/10.1016/S0378-1119(00)00428-5)
- Paradis, E., Claude, J., & Strimmer, K. (2004). APE: Analyses of phylogenetics and evolution in R language. *Bioinformatics*, *20*, 289–290. <https://doi.org/10.1093/bioinformatics/btg412>
- Perrière, G., & Thioulouse, J. (2002). Use and misuse of correspondence analysis in codon usage studies. *Nucleic Acids Research*, *30*, 4548–4555. <https://doi.org/10.1093/nar/gkf565>
- Pogoda, C. S., Keepers, K. G., Nadiadi, A. Y., Bailey, D. W., Lendemer, J. C., Tripp, E. A., & Kane, N. C. (2019). Genome streamlining via complete loss of introns has occurred multiple times in lichenized fungal mitochondria. *Ecology and Evolution*, *9*, 4245–4263. <https://doi.org/10.1002/ece3.5056>
- R Core Team (2018). *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing.
- Ranwez, V., Harispe, S., Delsuc, F., & Douzery, E. J. P. (2011). MACSE: Multiple alignment of coding SEquences accounting for frameshifts and stop codons. *PLoS One*, *6*, e22594. <https://doi.org/10.1371/journal.pone.0022594>
- Sánchez-Alonso, P., & Guzmán, P. (1998). Organization of chromosome ends in *Ustilago maydis*. RecQ-like helicase motifs at telomeric regions. *Genetics*, *148*, 1043–1054.
- Sanchez-Puerta, M. V., Abbona, C. C., Zhuo, S., Tepe, E. J., Bohs, L., Olmstead, R. G., & Palmer, J. D. (2011). Multiple recent horizontal transfers of the *cox1* intron in Solanaceae and extended co-conversion of flanking exons. *BMC Evolutionary Biology*, *11*, 277. <https://doi.org/10.1186/1471-2148-11-277>
- Sanchez-Puerta, M. V., Cho, Y., Mower, J. P., Alverson, A. J., & Palmer, J. D. (2008). Frequent, phylogenetically local horizontal transfer of the *cox1* group I Intron in flowering plant mitochondria. *Molecular Biology and Evolution*, *25*, 1762–1777. <https://doi.org/10.1093/molbev/msn129>
- Schirawski, J., Mannhaupt, G., Münch, K., Brefort, T., Schipper, K., Doehlemann, G., ... Kahmann, R. (2010). Pathogenicity determinants in smut fungi revealed by genome comparison. *Science*, *330*, 1546–1548. <https://doi.org/10.1126/science.1195330>
- Schuster, A., Lopez, J. V., Becking, L. E., Kelly, M., Pomponi, S. A., Wörheide, G., ... Cárdenas, P. (2017). Evolution of group I introns in Porifera: New evidence for intron mobility and implications for DNA barcoding. *BMC Evolutionary Biology*, *17*, 82. <https://doi.org/10.1186/s12862-017-0928-9>
- Schweizer, G., Münch, K., Mannhaupt, G., Schirawski, J., Kahmann, R., & Dutheil, J. Y. (2018). Positively selected effector genes and their contribution to virulence in the smut fungus *Sporisorium reilianum*. *Genome Biology and Evolution*, *10*, 629–645. <https://doi.org/10.1093/gbe/evy023>
- Sela, I., Ashkenazy, H., Katoh, K., & Pupko, T. (2015). GUIDANCE2: accurate detection of unreliable alignment regions accounting for the uncertainty of multiple parameters. *Nucleic Acids Res.*, *43*(Web Server issue), W7–W14. <https://doi.org/10.1093/nar/gkv318>
- Singh, D. K., Ghosh, A. K., Croteau, D. L., & Bohr, V. A. (2012). RecQ helicases in DNA double strand break repair and telomere maintenance.

- Mutation Research*, 736, 15–24. <https://doi.org/10.1016/j.mrfmmm.2011.06.002>
- Steinberg, G., & Perez-Martin, J. (2008). *Ustilago maydis*, a new fungal model system for cell biology. *Trends in Cell Biology*, 18, 61–67. <https://doi.org/10.1016/j.tcb.2007.11.008>
- Stoddard, B. L. (2005). Homing endonuclease structure and function. *Quarterly Reviews of Biophysics*, 38, 49–95. <https://doi.org/10.1017/S0033583505004063>
- Stone, C. L., Frederick, R. D., Tooley, P. W., Luster, D. G., Campos, B., Winegar, R. A., ... Blagden, T. (2018). Annotation and analysis of the mitochondrial genome of *Coniothyrium glycines*, causal agent of red leaf blotch of soybean, reveals an abundance of homing endonucleases. *PLoS One*, 13, e0207062.
- Sun, C. W., & Callis, J. (1993). Recent stable insertion of mitochondrial DNA into an *Arabidopsis polyubiquitin* gene by nonhomologous recombination. *The Plant Cell*, 5, 97–107.
- Thorsness, P. E., & Weber, E. R. (1996). Escape and migration of nucleic acids between chloroplasts, mitochondria, and the nucleus. *International Review of Cytology*, 165, 207–234.
- Valverde, M. E., Vandemark, G. J., Martínez, O., & Paredes-López, O. (2000). Genetic diversity of *Ustilago maydis* strains. *World Journal of Microbiology & Biotechnology*, 16, 49–55.
- Volff, J.-N. (2006). Turning junk into gold: Domestication of transposable elements and the creation of new genes in eukaryotes. *BioEssays*, 28, 913–922. <https://doi.org/10.1002/bies.20452>
- Vollmeister, E., Schipper, K., Baumann, S., Haag, C., Pohlmann, T., Stock, J., & Feldbrügge, M. (2012). Fungal development of the plant pathogen *Ustilago maydis*. *FEMS Microbiology Reviews*, 36, 59–77.
- Yang, Z. (2007). PAML 4: Phylogenetic analysis by maximum likelihood. *Molecular Biology and Evolution*, 24, 1586–1591. <https://doi.org/10.1093/molbev/msm088>
- Zhang, J., Nielsen, R., & Yang, Z. (2005). Evaluation of an improved branch-site likelihood method for detecting positive selection at the molecular level. *Molecular Biology and Evolution*, 22, 2472–2479. <https://doi.org/10.1093/molbev/msi237>

SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section.

How to cite this article: Dutheil JY, Münch K, Schotanus K, Stukenbrock EH, Kahmann R. The insertion of a mitochondrial selfish element into the nuclear genome and its consequences. *Ecol Evol*. 2020;10:11117–11132. <https://doi.org/10.1002/ece3.6749>