



# Coexistence of Multiple Endemic and Pandemic Lineages of the Rice Blast Pathogen

 Pierre Gladieux,<sup>a</sup>  Sébastien Ravel,<sup>a</sup> Adrien Rieux,<sup>b</sup> Sandrine Cros-Arteil,<sup>a</sup> Henri Adreit,<sup>a</sup> Joëlle Milazzo,<sup>a</sup> Maud Thierry,<sup>a</sup>  Elisabeth Fournier,<sup>a</sup>  Ryohei Terauchi,<sup>c</sup>  Didier Tharreau<sup>a</sup>

<sup>a</sup>UMR BGPI, Univ Montpellier, INRA, CIRAD, Montpellier SupAgro, Montpellier, France

<sup>b</sup>CIRAD, UMR PVBMT, St. Pierre de la Reunion, France

<sup>c</sup>Iwate Biotechnology Research Center, Kitakami, Iwate, Japan

**ABSTRACT** The rice blast fungus *Magnaporthe oryzae* (syn., *Pyricularia oryzae*) is both a threat to global food security and a model for plant pathology. Molecular pathologists need an accurate understanding of the origins and line of descent of *M. oryzae* populations in order to identify the genetic and functional bases of pathogen adaptation and to guide the development of more effective control strategies. We used a whole-genome sequence analysis of samples from different times and places to infer details about the genetic makeup of *M. oryzae* from a global collection of isolates. Analyses of population structure identified six lineages within *M. oryzae*, including two pandemic on japonica and indica rice, respectively, and four lineages with more restricted distributions. Tip-dating calibration indicated that *M. oryzae* lineages separated about a millennium ago, long after the initial domestication of rice. The major lineage endemic to continental Southeast Asia displayed signatures of sexual recombination and evidence of DNA acquisition from multiple lineages. Tests for weak natural selection revealed that the pandemic spread of clonal lineages entailed an evolutionary “cost,” in terms of the accumulation of deleterious mutations. Our findings reveal the coexistence of multiple endemic and pandemic lineages with contrasting population and genetic characteristics within a widely distributed pathogen.

**IMPORTANCE** The rice blast fungus *Magnaporthe oryzae* (syn., *Pyricularia oryzae*) is a textbook example of a rapidly adapting pathogen, and it is responsible for one of the most damaging diseases of rice. Improvements in our understanding of *Magnaporthe oryzae*'s diversity and evolution are required to guide the development of more effective control strategies. We used genome sequencing data for samples from around the world to infer the evolutionary history of *M. oryzae*. We found that *M. oryzae* diversified about 1,000 years ago, separating into six main lineages: two pandemic on japonica and indica rice, respectively, and four with more restricted distributions. We also found that a lineage endemic to continental Southeast Asia displayed signatures of sexual recombination and the acquisition of genetic material from multiple lineages. This work provides a population-level genomic framework for defining molecular markers for the control of rice blast and investigations of the molecular basis of differences in pathogenicity between *M. oryzae* lineages.

**KEYWORDS** clonality, deleterious mutations, indica rice, introgression, japonica rice, population genomics, population structure, recombination, rice blast, tip-dating calibration

Fungal plant pathogens provide many examples of geographically widespread, often clonal, lineages capable of adapting rapidly to anthropogenic changes, such as the use of new fungicides or resistant varieties, despite extremely low levels of population

Received 3 October 2017 Accepted 26 February 2018 Published 3 April 2018

**Citation** Gladieux P, Ravel S, Rieux A, Cros-Arteil S, Adreit H, Milazzo J, Thierry M, Fournier E, Terauchi R, Tharreau D. 2018. Coexistence of multiple endemic and pandemic lineages of the rice blast pathogen. mBio 9:e01806-17. <https://doi.org/10.1128/mBio.01806-17>.

**Editor** David Guttman, University of Toronto

**Copyright** © 2018 Gladieux et al. This is an open-access article distributed under the terms of the [Creative Commons Attribution 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

Address correspondence to Pierre Gladieux, [pierre.gladieux@inra.fr](mailto:pierre.gladieux@inra.fr).

genetic diversity (1, 2). An accurate characterization of the population biology and evolutionary history of these organisms is crucial to an understanding of the factors underlying their emergence and spread and to provide new, powerful, and enduring solutions to control these factors. Knowledge of the origins and lines of descent connecting extant pathogen populations provides insight into the pace and mode of disease emergence and subsequent dispersal (2, 3). By inferring the history and structure of pathogen populations, we can also identify disease reservoirs and improve our understanding of the transmissibility and longevity of populations (4, 5). Finally, quantification of the amount and distribution of genetic variation across space and time provides a population-level genomic framework for defining molecular markers for pathogen control and for investigations of the molecular basis of differences in phenotype and fitness between divergent pathogen lineages.

Rice blast is one of the most damaging rice diseases worldwide (6–8). It is caused by the ascomycete fungus *Magnaporthe oryzae* (syn., *Pyricularia oryzae*), which has become a model for plant pathology in parallel with the development of rice as a model crop species (7, 9–11). The rice-infecting lineage of *M. oryzae* coexists with multiple host-specialized and genetically divergent lineages that infect other cereals and grasses (12–14). The lineage infecting foxtail millet (*Setaria italica*, referred to hereafter as *Setaria*) is the closest relative of the rice-infecting lineage, and rice blast was thus thought to have emerged following a host shift from *Setaria* about 2,500 to 7,500 years ago (15), at a time when *Setaria* was the preferred staple in East Asia (16, 17). *Magnaporthe oryzae* infects the two major subspecies of rice, *Oryza sativa* subsp. *indica* and *Oryza sativa* subsp. *japonica* (referred to here as *indica* and *japonica*, respectively). Population genomics studies have provided support for a model in which *de novo* domestication occurred only once, to generate the *japonica* lineage, which subsequently diverged into temperate and tropical *japonica*, with introgressive hybridization from *japonica* leading to domesticated *indica* (18–20). Using microsatellite markers, Saleh et al. (21) identified multiple endemic and pandemic genetic pools of rice-infecting strains, but they were unable to resolve the evolutionary relationships between them. Rice blast has proved able to adapt rapidly to varietal resistance and is thus a dynamic threat to such resistance in rice agrosystems (22). This ability to adapt is surprising given the low level of diversity in *M. oryzae* and its infertility or asexual mode of reproduction in most rice-growing areas (22, 23). This pathogen may thus be particularly exposed to the “cost of pestification” (by analogy with the cost of domestication [24–27]), according to which the combination of a small effective population size, strong selection on pestification genes, and a lack of recombination lead to the accumulation of deleterious mutations (28). Potential limitations to adaptation could be counterbalanced by boom-and-bust cycles in *M. oryzae*, with adaptation occurring during the boom phases, when the short-term effective population size is large (2, 29). Adaptive mutations may also be introduced by cryptic genetic exchanges with conspecifics or heterospecifics (30–33), but these mechanisms remain to be investigated in natural populations of *M. oryzae* (34). An accurate understanding of the population genetics of successful clonal fungal pathogens, such as *M. oryzae*, can provide important insights into the genomic and eco-evolutionary processes underlying pathogen emergence and adaptation to anthropogenic changes.

We used pathogenicity data and whole-genome resequencing data for *M. oryzae* samples distributed over time and space to address the following questions. What population structure does *M. oryzae* display? Does this species consist of relatively ancient or recent clonal lineages? What is the history of temperate *japonica*, tropical *japonica*, and *indica japonica* rice colonization by *M. oryzae*? Do *M. oryzae* lineages display differences in pathogenicity toward rice subspecies? Can we identify genetic exchanges between rice-infecting lineages and the genomic regions that have been exchanged? Is there evidence for a cost of pestification in terms of the accumulation of deleterious mutations?

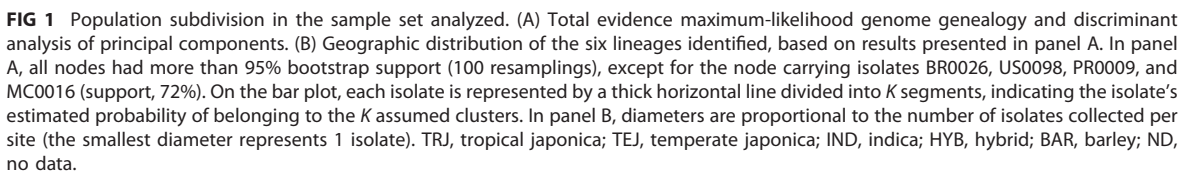
## RESULTS

**Genome sequencing and SNP calling.** We elucidated the emergence, diversification, and spread of *M. oryzae* in rice agrosystems by studying genome-wide variation across geographically widespread samples. We used 25 and 18 genomes sequenced by Illumina single-end and paired-end read technologies, respectively, with 7 published genome sequences obtained via Solexa and mate-pair titanium methods (10, 12). We thus had a total of 50 genomes available for analysis (see Table S1 in the supplemental material). Forty-five of the isolates concerned originated from cultivated rice (*Oryza sativa*), four from cultivated barley (*Hordeum vulgare*), and one from foxtail millet (*Setaria italica*). The sample set included multiple samples from geographically separated areas (North and South America, South, Southeast, and East Asia, sub-Saharan Africa, Europe, and the Mediterranean), and the reference laboratory strain 70-15 and its parent GY11 were from French Guiana. Nine samples were collected from tropical japonica rice, 7 from temperate japonica, 15 from indica, and 3 from hybrid elite varieties. Sequencing reads were mapped onto the 41.1-Mb reference genome of strain 70-15. Mean sequencing depth ranged from 5× to 64× for genomes sequenced with single-end reads and from 5× to 10× for genomes sequenced with paired-end reads (Table S1). Single-nucleotide polymorphism (SNP) calling identified 182,804 biallelic SNPs distributed over seven chromosomes. The data set consisted of 95,925 SNPs, excluding the *Setaria*-infecting lineage, 61,765 of which had less than 30% missing data and 16,370 of which had no missing data.

**Population subdivision, genealogical relationships, and levels of genetic variation.** We used a multivariate analysis of population subdivision method, rather than model-based clustering algorithms, because multivariate methods require no assumptions about outcrossing, random mating, or linkage equilibrium within clusters, and previous studies have shown that, in many populations, *M. oryzae* has lost its sexual recombination capacity (references 21 to 23 and references therein). We used a discriminant analysis of principal components (DAPC) to determine the number of lineages represented in our data set. When we progressively increased the number of clusters ( $K$ ) from 2 to 5, we identified the four lineages previously described by Saleh et al. in Asia (21) on the basis of microsatellite data, and we also identified a cluster of three strains collected from the Yunnan and Hunan provinces of China (Fig. 1). Further increases in  $K$  led to the subdivision of this Yunnan-Hunan cluster. Barley-infecting isolates clustered within rice-infecting lineage 1, which confirmed findings of previous phylogenetic studies (12, 13). Barley is “universally susceptible” to rice-infecting isolates, at least under laboratory conditions. However, the barley isolates included in this study were collected in Thailand, and no major blast epidemic has since been reported on this host in this area, indicating that barley is a minor host for rice-infecting populations.

We investigated whether the clusters observed at  $K$  values of  $>4$  in the DAPC represented new independent lineages or subdivisions of the main clusters by using RAxML to infer a genome genealogy (35). We based the analysis on a data set combining the full set of SNPs and monomorphic sites, rather than just SNPs, to increase topological and branch length accuracy (36). The total evidence genealogy revealed the existence of four lineages, corresponding to lineages 1 to 4 described by Saleh et al. (21), and two new lineages (named lineages 5 and 6) corresponding to the three-individual cluster observed at  $K = 5$  in the DAPC (Fig. 1). With the 41-Mb data set, including missing data, the most basal divergence within the rice-infecting lineage was that between lineage 1 and the other five lineages (Fig. 1). If positions with missing data were excluded (15 Mb), the most basal divergence was that between a group composed of lineages 1, 2, and 6 and a group composed of lineages 3, 4, and 5 (data not shown).

Absolute divergence ( $d_{xy}$ ) between pairs of lineages was on the order of  $10^{-4}$  differences per base pair and was highest in comparisons with lineage 6 (Table S2). Nucleotide diversity within lineages was an order of magnitude lower than divergence



**TABLE 1** Summary of population genomic variations in nonoverlapping 100-kb windows<sup>a</sup>

Lineage	<i>n</i>	<i>S</i>	<i>K</i>	<i>H<sub>e</sub></i>	$\theta_w$	$\pi$	<i>D</i>
1	10	57.6	3.6	0.31	2.25E-04	2.11E-04	-0.558
2	14	19.7	7.2	0.20	6.94E-05	4.92E-05	-1.454
3	16	21.5	7.9	0.17	7.24E-05	4.53E-05	-1.718
4	6	10.5	4.3	0.38	5.15E-05	4.52E-05	-0.824

<sup>a</sup>Lineages 5 and 6 were not included in calculations because the sample sizes for these lineages were too small (*n* = 2 and *n* = 1, respectively). *n*, sample size;  $\theta_w$ , Watterson's  $\theta$  per base pair;  $\pi$ , nucleotide diversity per base pair; *H<sub>e</sub>*, haplotype diversity; *K*, number of haplotypes; *D*, Tajima's neutrality statistic.

in lineages 2 to 4 ( $\theta_w$  per site, 5.2e-5 to 7.2e-5;  $\pi$  per site, 4.5e-5 to 4.9e-5) and was highest in lineage 1 ( $\theta_w$  per site, 2.3e-4;  $\pi$  per site, 2.1e-4) (Table 1). Tajima's *D* was negative in all lineages, indicating an excess of low-frequency polymorphisms, and values were closer to zero in lineages 1 and 4 (*D* = -0.56 and -0.82, respectively) than in lineages 2 and 3 (*D* = -1.45 and -1.72, respectively). The same differences in levels of variability across lineages, and individual summary statistics of the same order of magnitude, were observed if missing data were excluded from computations.

**Footprints of natural selection and the cost of pestification.** We tested for standard neutral molecular evolution by using the McDonald-Kreitman method, based on genome-wide patterns of synonymous and nonsynonymous variations (Table 2). The null hypothesis could be rejected for all four lineages (*P* < 0.0001). The neutrality index, which quantifies the direction and degree of departure from neutrality, was greater than 1, indicating an excess of amino acid polymorphisms. This pattern suggests that lineages 1 to 4 accumulated slightly deleterious mutations during divergence from the *Setaria*-infecting lineage. Under near-neutrality, the ratio of nonsynonymous to synonymous nucleotide diversity ( $\pi_N/\pi_S$ ) provides an estimate of the proportion of effectively neutral mutations that are strongly dependent on the effective population size, *N<sub>e</sub>* (37). The  $\pi_N/\pi_S$  ratio ranged from 0.43 in lineage 1 to 0.61 in lineage 4 and was intermediate in lineages 2 and 3 ( $\pi_N/\pi_S$  = 0.49), and the ratio of nonsense (i.e., premature stop codons) to sense nonsynonymous mutations (*P<sub>nonsense</sub>*/*P<sub>sense</sub>*) followed the same pattern. Overall, the  $\pi_N/\pi_S$  and *P<sub>nonsense</sub>*/*P<sub>sense</sub>* ratios obtained suggest a higher proportion of slightly deleterious mutations segregating in lineage 4 and, to a lesser extent, in lineages 2 and 3, than in lineage 1. Assuming identical mutation rates, we can estimate that the long-term population size of lineage 1 ( $\pi_S$  = 0.00018/bp) was 2.5 to 3 times greater than that of the other lineages, consistent with the effect of *N<sub>e</sub>* on the efficacy of negative selection predicted under near-neutrality.

**Distribution and reproductive biology of *M. oryzae* lineages.** The strains of lineages 1 and 2 originated from rain-fed upland rice, including rice grown in experimental fields. Lineage 2 was exclusively associated with tropical and temperate japonica, whereas lineage 1 was sampled from barley, tropical japonica, and hybrid rice varieties (Fig. 1; Table S1). Lineage 1 was restricted to continental Southeast Asia (Laos, Thailand, Yunnan). The reference laboratory strain GY-11 (also referred to as Guy11) was

**TABLE 2** Results of McDonald-Kreitman tests based on genome-wide patterns of synonymous and nonsynonymous variation and measurements of the genome-wide intensity of purifying selection<sup>a</sup>

Lineage	$\pi_N/\pi_S$	<i>P<sub>nonsense</sub></i> / <i>P<sub>sense</sub></i>	<i>P<sub>n</sub></i> / <i>P<sub>s</sub></i>	<i>D<sub>n</sub></i> / <i>D<sub>s</sub></i>	NI
1	0.43 (0.00041/0.00018)	0.011 (49/4,244)	1.23 (4,293/3,492)	0.70 (16,444/23,656)	1.77*
2	0.49 (0.00012/0.00006)	0.022 (36/1,622)	1.52 (1,658/1,088)	0.72 (15,565/21,745)	2.13*
3	0.49 (0.00015/0.00007)	0.018 (32/1814)	1.17 (1,846/1,578)	0.97 (14,789/15,293)	1.21*
4	0.61 (0.00012/0.00007)	0.034 (31/914)	1.59 (945/593)	0.72 (15,302/21,347)	2.22*

<sup>a</sup>Divergence was measured against predicted gene sequences of the *Setaria*-infecting *Magnaporthe oryzae* isolate US71.  $\pi_N/\pi_S$  is the ratio of nonsynonymous to synonymous nucleotide diversity. Under near-neutrality,  $\pi_N/\pi_S$  provides an estimate of the proportion of effectively neutral mutations strongly dependent on effective population size, *N<sub>e</sub>*.  $\pi_S$  is a proxy for *N<sub>e</sub>*. *P<sub>nonsense</sub>*/*P<sub>sense</sub>* is the number of nonsynonymous nonsense mutations (e.g., a "premature" stop codon) divided by the number of nonsynonymous sense mutations. The neutrality index (NI) = (*P<sub>n</sub>*/*P<sub>s</sub>*)/(*D<sub>n</sub>*/*D<sub>s</sub>*) and determines the direction and degree of departure from neutrality; \*, *P* < 0.0001, chi-square test of independence. NI is equal to 1 if nonsynonymous mutations are neutral or strongly deleterious. NI is <1 when amino acid substitutions have occurred and implies that advantageous mutations have become fixed. NI is >1 when there is an excess of amino acid polymorphisms, as expected in a context of slightly deleterious mutations.



collected in French Guiana, from fields cultivated by Hmong refugees who fled Laos in the 1970s. Lineage 2 was pandemic and included all the European samples.

Lineage 3 and 4 samples originated from irrigated or rain-fed upland/lowland rice. They were mostly associated with indica rice, with two samples collected from hybrid varieties and one collected from tropical japonica (Fig. 1; Table S1). Lineage 3 was pandemic and was found in all sub-Saharan Africa samples, whereas lineage 4 was found on the Indian subcontinent, in Zhejiang (China), and the United States. Lineages 5 and 6 were collected from indica and tropical japonica varieties of rain-fed upland rice in Yunnan and Hunan, China, respectively.

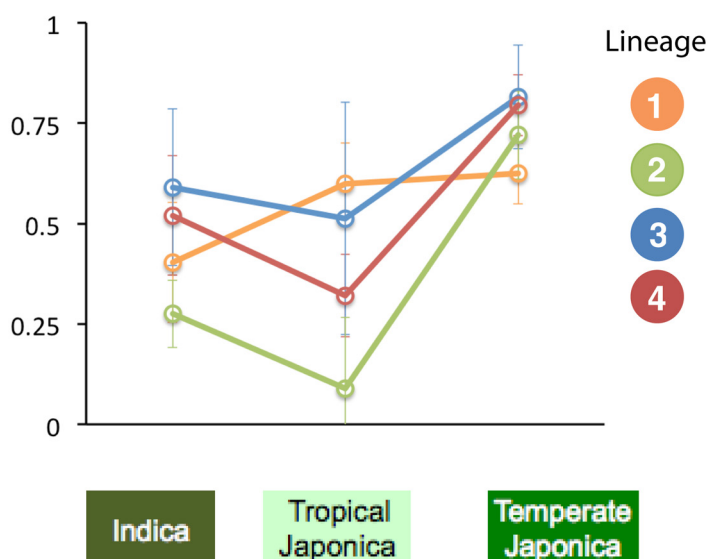
Lineages 2, 3, and 4 displayed low rates of female fertility (20%, 0%, and 0%, respectively) and a significant imbalance in mating type ratio (frequency of Mat-1, 100%, 14.3%, and 100%, respectively; chi-square test,  $P < 0.001$ ), whereas lineage 1 had a female fertility rate of 88.9% and a nonsignificant imbalance in mating type ratio (frequency of Mat-1, 33.3%; chi-square test,  $P = 0.083$ ). Lineage 5 was Mat-1, and only one of the two strains was female fertile (no data for lineage 6).

**Pathogen compatibility range.** Gallet et al. (38) analyzed the range of compatibility, in terms of the qualitative success of infection, between 31 *M. oryzae* isolates and 57 rice genotypes. Analyses of variance revealed a pattern of host-pathogen compatibility strongly structured by the host of origin of the isolates (i.e., the rice subspecies from which samples were collected). We investigated whether the compatibility between rice hosts and *M. oryzae* isolates was also structured by the lineage of origin of the isolates, by supplementing the data set published by Gallet et al. (38) with pathotyping data for 27 isolates. We added microsatellite data to the SNP data, to overcome the absence of sequence data for 28 isolates, and we used clustering methods to confidently assign 46 of the 58 isolates with pathotyping data to identified lineages (no isolates could be assigned to lineage 5 or 6 [see Materials and Methods]). The final pathogenicity data set included 46 isolates from lineages 1 to 4, inoculated onto 38 tropical japonica, temperate japonica, and indica varieties and 19 differential varieties with known resistance genes (Table S3).

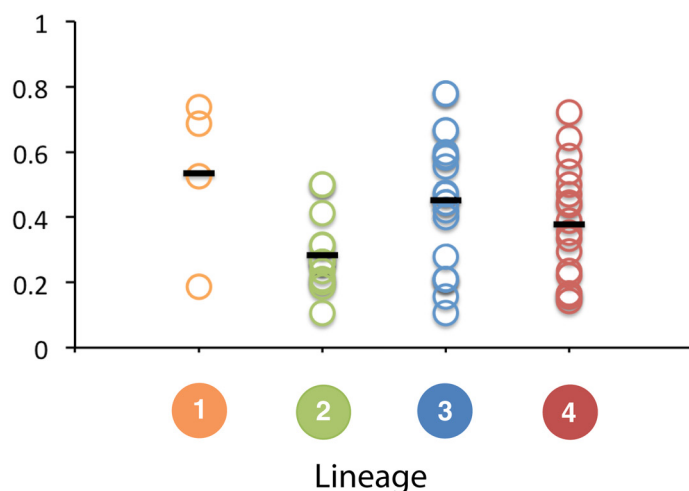
Infection success (binary response) was analyzed with a generalized linear model. An analysis of the proportion of compatible interactions revealed significant effects of rice subspecies, pathogen lineage, and the interaction between them (Table S4). The lineage effect could be explained by lineage 2 having a lower infection frequency than lineage 1 (comparison of lineages 1 and 2:  $z = -2.779$ ,  $P = 0.005$ ) and by lineage 3 having a higher infection frequency than lineage 1 (comparison of lineages 3 and 1:  $z = 2.683$ ,  $P = 0.007$ ), whereas the infection frequency of lineage 4 was not significantly different from that of lineage 1 (comparison of lineages 4 and 1:  $z = 1.121$ ,  $P = 0.262$ ). The rice subspecies effect could be attributed to tropical japonica varieties having a wider compatibility range than indica varieties (comparison of tropical japonica and indica:  $z = 1.793$ ,  $P = 0.073$ ) and temperate japonica having a wider compatibility range than indica varieties (comparison of temperate japonica and indica:  $z = 1.830$ ,  $P = 0.067$ ). The significant interaction between rice subspecies and pathogen lineage indicates that the effect of the lineage of origin of the isolate on the proportion of compatible interactions differed between the three rice subspecies. This interaction effect can be attributed to pathogen specialization on indica and tropical japonica, with lineage 1 (mostly originating from tropical japonica or from areas in which tropical japonica is grown) infecting tropical japonica varieties more frequently than indica varieties, lineage 2 (the lineage sampled from temperate japonica) infecting temperate japonica varieties more frequently than other varieties, lineages 3 and 4 (mostly originating from indica varieties) infecting indica varieties more frequently than tropical japonica varieties, and all four lineages infecting temperate japonica varieties at relatively high frequencies (Fig. 2A; Table S4).

Major resistance (R) genes can be a major determinant of pathogen host range, and they promote divergence between pathogen lineages by exerting strong divergent selection on a limited number of pathogenicity-related genes (39–41). We investigated

## (A) Proportion of compatible interactions



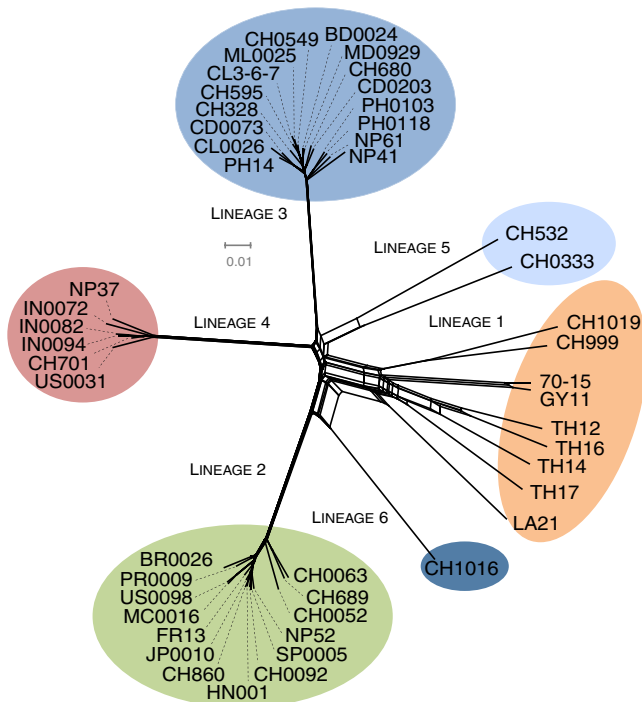
## (B) Proportion of R genes overcome



**FIG 2** Proportion of compatible interactions between 46 isolates from lineages 1 to 4 of *M. oryzae* and 38 varieties representing three rice subspecies (A) and the proportion of R genes overcome by 36 isolates from lineages 1 to 4 of *M. oryzae* used to inoculate 19 differential lines of rice (B).

the possible role of major resistance genes in the observed differences for compatibility between rice subspecies and pathogen lineages by challenging 19 differential varieties with the 46 isolates assigned to lineages 1 to 4. An analysis of the number of R genes overcome revealed a significant effect of pathogen lineage (Table S5). This effect was driven mostly by lineage 2, which overcame fewer R genes than the other lineages (Fig. 2B; Table S5).

**Recombination within and between lineages.** We visualized evolutionary relationships while taking into account the possibility of recombination within or between lineages by using the phylogenetic network approach Neighbor-Net, as implemented in Splitstree 4.13 (42). Neighbor-Net is an agglomerative method that generates planar split graph representations. A split is a partitioning of the data set, and a collection of splits is considered compatible if they fall within the set of splits of a tree. Gene genealogies represent compatible collections of splits, whereas Neighbor-Net can be



**FIG 3** Neighbor-Net networks showing relationships between haplotypes identified on the basis of the full set of 16,370 SNPs without missing data in the whole sample set (A), in lineage 1 (B), in lineage 2 (C), in lineage 3 (D), and in lineage 4 (E).

used to visualize conflicting phylogenetic signals, represented by network reticulation, through a condition weaker than compatibility. The Neighbor-Net network inferred from the set of 16,370 SNPs without missing data presented a non-tree-like structure of the inner connections between lineages, consistent with genetic exchanges between unrelated isolates or incomplete lineage sorting (Fig. 3). Greater network reticulation was observed between lineages 1, 5, or 6 and the other lineages than between these other lineages themselves. Lineages 2 to 4 had long interior branches and star-like topologies, consistent with long-term clonality.

We evaluated the amount of recombination within lineages by estimating the population recombination parameter ( $\rho = 2 N_e r$ ) and testing for the presence of recombination with a likelihood permutation test implemented in the Pairwise program in LDHAT. Recombination analyses confirmed the heterogeneity between lineages of the contribution of recombination to genomic variation, with recombination rates averaged across chromosomes of more than 2 to 3 orders of magnitude higher in lineage 1 (10.57 crossovers/Mbp/generation) than in other lineages (lineage 2, 0.28; lineage 3, 0.01; lineage 4, 0.33 crossovers/Mbp/generation) (Table 3). SplitsTree analy-

**TABLE 3** Estimates of the population recombination rate ( $\rho$ ), tests of recombination based on homoplasy and linkage disequilibrium, and the proportion of homoplastic SNPs

	$\rho$ (no. of crossovers/Mbp/generation) on chromosome <sup>a</sup>								% homoplastic SNPs	phi test P value
Lineage	1	2	3	4	5	6	7	Mean		
1	8.6*	3.8*	15.1*	1.4	8.5*	10.6*	13.5*	10.57	34.56	0.0000
2	0.0	0.2	0.2*	0.3	0.6	0.5	0.3	0.28	0.09	0.0944
3	0.4*	0.2*	0.0	0.0	0.0	0.0	0.0	0.01	0.47	0.0535
4	0.2	0.2*	0.3	0.4	0.4	0.4	0.4	0.33	0.40	0.0014

<sup>a</sup>\*,  $P < 0.05$ . The phi test assesses pairwise homoplasy. The null hypothesis of no recombination was tested, with the phi test and for  $\rho$ , using random permutations of the positions of the SNPs based on the expectation that sites are exchangeable if there is no recombination. For the  $\rho$  test, significance was determined from the distribution of maximum composite likelihood values calculated from permuted data.



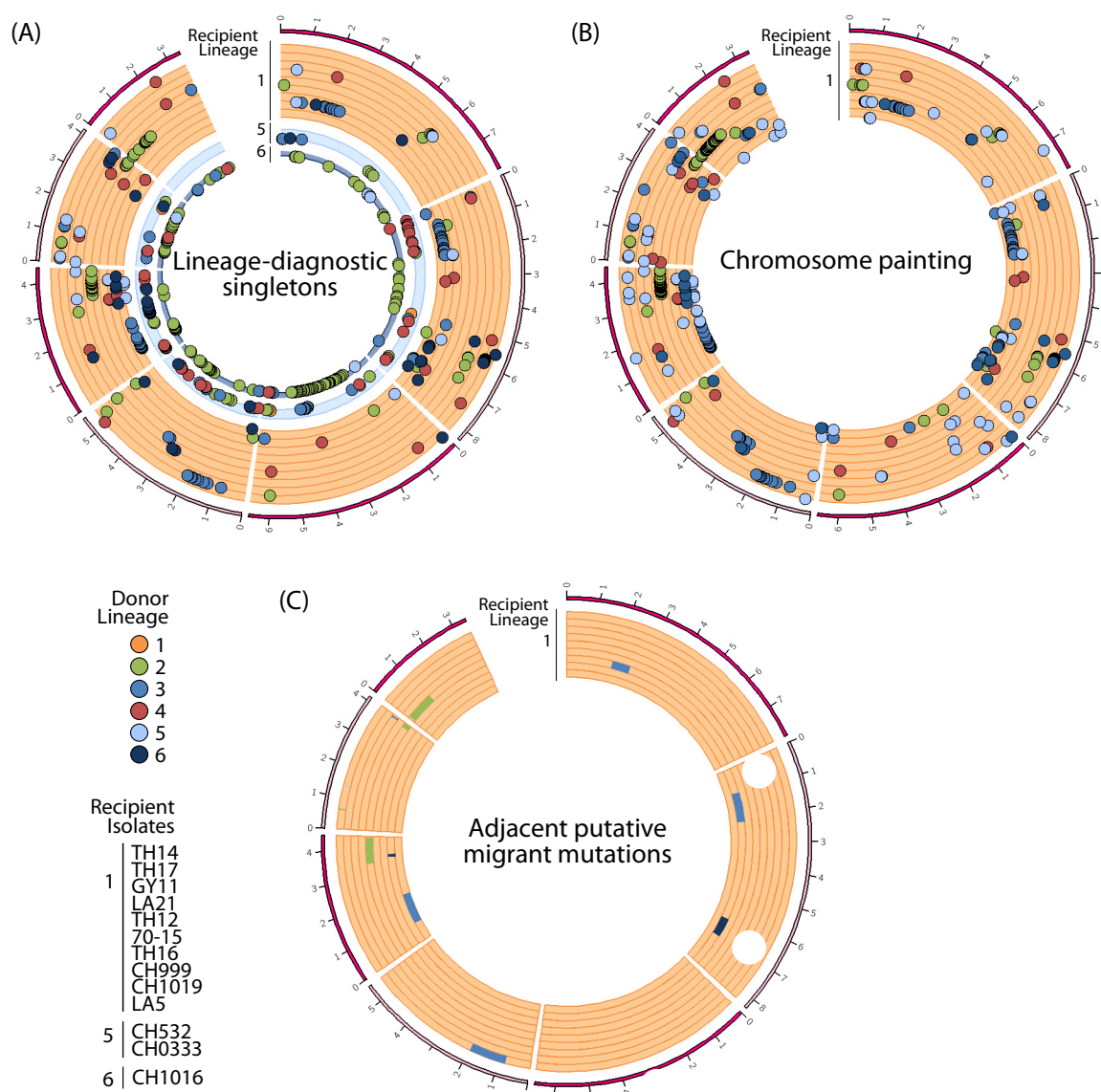
ses, producing the reticulations within each lineage and testing for recombination with the phi test, were consistent with this pattern (Table 3 and Fig. 3; Fig. S1). The null hypothesis of no recombination was rejected only for lineages 1 and 4 (43).

Differences in recombination-based variation between lineages were confirmed by analyses of homoplasy (Table 3). Homoplastic sites display sequence similarities that are not inherited from a common ancestor; instead, they result from independent events in different branches. Homoplasy can result from recurrent mutations or recombination, and the contribution of recombination to homoplasy is expected to predominate in outbreeding populations. Homoplastic sites were identified by mapping mutations onto the total evidence genome genealogy with the Trace All Characters function of Mesquite (44), applying ancestral reconstruction under the maximum parsimony optimality criterion. The resulting matrix of ancestral states for all nodes was then processed with a python script to determine the number of mutations that had occurred at each site within each lineage, counting sites displaying multiple substitutions across the tree as homoplastic. Only 0.09%, 0.47%, and 0.40% of the SNPs were homoplastic in lineages 2, 3, and 4, respectively, versus 34.6% in lineage 1 (lineages 5 and 6 were not tested due to the small sample sizes). The very small numbers of homoplastic sites in lineages 2, 3, and 4 suggested that these lineages are largely clonal, whereas the high level of homoplasy detected in lineage 1 is consistent with repeated recombination events between strains of this lineage.

We assessed the genomic impact of recombination by analyzing patterns of linkage disequilibrium (LD), i.e., the tendency of different alleles to occur together in a non-random manner. For lineage 1 ( $S = 13,000$  SNPs), LD decayed smoothly with physical distance, reaching half its maximum value at about 10 kb, whereas for lineages 2, 3, and 4 ( $S = 3,700$ , 3,200, and 2,700 SNPs), no LD decay pattern was observed (Fig. S2). These analyses also revealed that background LD levels were no higher in lineages 2, 3, or 4, which appeared to be largely clonal, than in lineage 1. However, both simulation work and empirical data have shown that population history, including bottlenecks and admixtures, strongly affects the background level of LD in a population (45).

**Genome scan for genetic exchanges between lineages.** We scanned the genomes for the exchange of mutations between lineages, using a method based on lineage-diagnostic SNPs and a probabilistic method of “chromosome painting” (Fig. 4). In the lineage-diagnostic SNP approach, each isolate is removed from the data set in turn to identify SNPs specific to a particular lineage (i.e., biallelic sites displaying a mutation specific to a given lineage). Each focal isolate is then added back to the data set and scanned for the presence of lineage-diagnostic SNPs identified in lineages other than its lineage of origin. Using this approach, we identified 515 lineage-diagnostic singletons with 276, 96, and 140 singletons in lineages 1, 5, and 6, respectively, and only 1 singleton in each of lineages 2, 3, and 4. Putatively migrant singletons were assigned to all other lineages for lineages 1 and 5 and to all other lineages except lineage 1 for lineage 6 (Table S6). Chromosome painting is a probabilistic method for reconstructing the chromosomes of each individual sample as a combination of all other homologous sequences. We identified the migrant mutations present in each isolate, with these mutations being defined as those having a probability greater than 90% of resulting from being copied from a lineage other than the lineage of origin of the focal isolate. This method uses population data from recipient populations only, and we were therefore able to include only lineages 1 to 4 in the analysis. Chromosome painting identified 464 migrant mutations, all of which segregated in lineage 1. Putative migrant mutations were assigned to all five of the other lineages (92.8 mutations per lineage, on average), with lineage 2 making the largest contribution (165 mutations) and lineage 4 the smallest contribution (39 mutations).

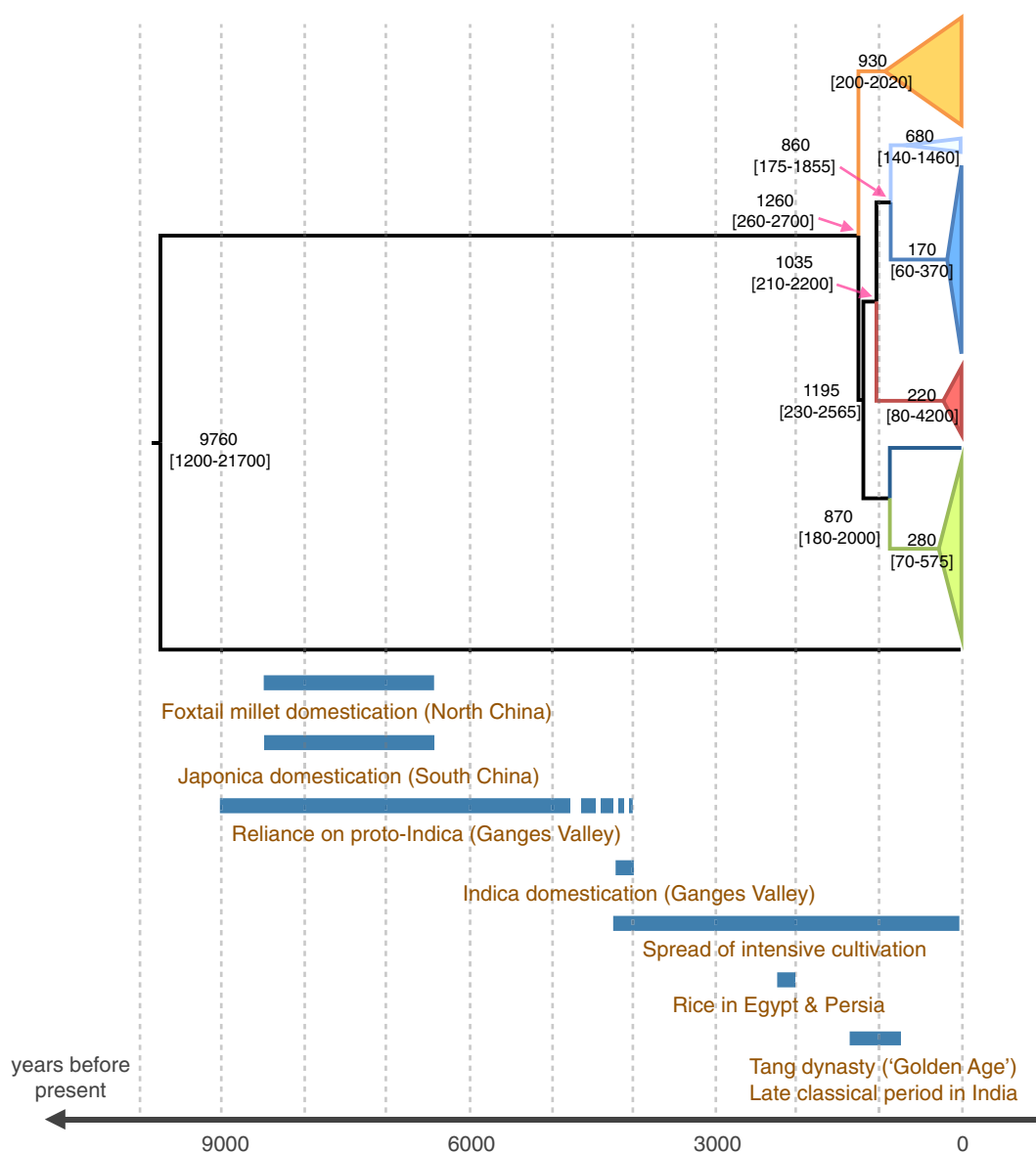
The sets of putative migrant mutations identified by the two methods matched different sets of genes enriched in NOD-like receptor (NLR) (46), HET domain (47), or the GO term lipid catabolic process (Table S6). However, the presence of false positives due to the random sorting of ancestral polymorphisms in lineage 1 and other lineages



**FIG 4** Genomic distribution of candidate immigrant mutations in lineages 1, 5, and 6. (A) Lineage-diagnostic mutations segregating as singletons in other lineages. (B) Lineage 1 mutations for which the most probable donor was lineage 2, 3, or 4 in probabilistic chromosome painting analysis. Lineages 5 and 6 ( $n = 2$  and  $n = 1$ , respectively) could not be included as recipient populations in the chromosome painting analysis due to their small sample sizes. No candidate immigrant mutations were identified in lineages 2, 3, or 4. (C) Genomic regions corresponding to a series of adjacent putative migrant mutations identified with lineage-diagnostic singletons in lineage 1. Chromosomes 1 to 7 appear in clockwise order, with ticks at megabase intervals.

cannot be excluded. We minimized the impact of the retention of ancestral mutations by reasoning that series of adjacent mutations are more likely to represent genuine gene exchange events. We identified all the genomic regions defined by three adjacent putative migrant mutations originating from the same donor lineage. We searched for such mutations among the set of putative migrant mutations identified by the two methods. We identified 12 such regions in total, corresponding to 1,917 genes. Functional enrichment tests for each recipient isolate revealed enrichment of genes for the GO term pathogenesis for isolate CH999, the GO term phosphatidylinositol biosynthetic process for isolate TH17, and the GO term telomere maintenance for isolate CH1019 (Table S6).

**Molecular dating.** We investigated the timing of rice blast emergence and diversification by performing Bayesian phylogenetic analyses with Beast. Isolates were collected from 1967 to 2009 (Table S1), making it possible to use a tip-based calibration



**FIG 5** Tip-calibrated genealogy inferred by maximum-likelihood phylogenetic inference using Beast 1.8.2, based on single-nucleotide variations in 50 *M. oryzae* genomes. Approximate historical periods are shown for context.

approach to estimate evolutionary rates and ancestral divergence times together. We analyzed the linear regression of sample age against root-to-tip distance (i.e., the number of substitutions separating each sample from the hypothetical ancestor at the root of the tree). The temporal signal obtained in this analysis was strong enough for thorough tip-dating inferences (Fig. S3) (48). We therefore used tip dating to estimate the rate at which mutations accumulate (i.e., the substitution rate) and the age of every node in the tree, including the root (i.e., time to the most recent common ancestor), simultaneously. At the scale of the genome, the mean substitution rate was estimated at  $1.98 \times 10^{-8}$  substitutions/site/year (Fig. S3). The six rice-infecting lineages were estimated to have diversified ~900 to ~1,300 years ago (95% highest posterior density [HPD], 175 to 2,700 years ago) (Fig. 5). Bootstrap node support was strong, and similar node age estimates were obtained when the recombining lineage 1 and the potentially recombining lineages 5 and 6 (data not shown) were excluded, indicating the limited effect of recombination on our inferences. We also inferred that the ancestor of rice-infecting and *Setaria*-infecting lineages lived ~9,800 years ago. However, the

credibility intervals were relatively large (95% HPD, 1,200 to 22,000 years ago), covering the period from japonica rice domestication and *Setaria* domestication to the last glacial maximum and overlapping with previous estimates suggesting that the rice- and *Setaria*-infecting lineages diverged shortly after rice domestication, or even during the period of rice domestication (range of point estimates in reference 15, 2,500 to 7,300 years ago).

## DISCUSSION

We performed a whole-genome sequence analysis of 50 isolates with different temporal and spatial distributions in order to elucidate the emergence, diversification, and spread of *M. oryzae* as a rapidly evolving pathogen with a devastating impact on rice agrosystems. Analyses of population subdivision confirmed the four lineages previously identified by Saleh et al. (21). Previous analyses of microsatellite data were unable to resolve the genealogical relationships between clusters or to capture the phylogenetic depth of population subdivision within *M. oryzae*. In contrast, our population genomic analyses of resequencing data revealed weak divergence between clusters (absolute divergence [ $d_{xy}$ ] on the order of  $10^{-4}$  differences per base pair), consistent with recent diversification. Phylogenetic analyses using sampling dates for calibration confirmed the recent origin of the six lineages, with estimates of divergence times ranging from ~900 to ~1,300 years ago (95% credible intervals, 175 to 2,700 years ago). Lineage 1 (which includes the reference strains GY11 and 70-15) was found in mainland Southeast Asia and originates from barley, tropical japonica, or undetermined varieties. All isolates from lineages 1, 5, and 6 were collected in rain-fed upland agrosystems typical of japonica rice cultivation, and pathogenicity test results were consistent with the local adaptation of lineage 1 to tropical japonica rice. Lineage 2 was pandemic in irrigated fields of temperate japonica rice outside Asia, and cross-inoculation experiments revealed specialization on this host and an ability to overcome fewer R genes, on average, than other lineages. Lineages 3 and 4 were associated with indica. Lineage 3 is pandemic, and cross-inoculation indicated local adaptation to this host, relative to tropical japonica, although lineages 3 and 4 had relatively wide compatibility ranges, consistent with generalism. One possible explanation for the wide compatibility range of temperate japonica varieties and the narrow compatibility range of lineage 2 is that temperate japonica varieties have smaller repertoires of R genes, as resistance to blast is of less concern to breeders growing rice under temperate irrigated conditions, which are less conducive to epidemics (38).

The continental Southeast Asian lineage was the most basal in total evidence genome genealogies, reflecting a pathway of domesticated Asian rice evolution (16, 18) in which the *de novo* domestication of rice occurred only once, in japonica. However, the diversification of *M. oryzae* into multiple rice-infecting lineages (point estimates ranging from ~900 to ~1,300 years ago) appears to be much more recent than the *de novo* domestication of rice (8,500 to 6,500 years ago [16, 49, 50]), the spread of rice cultivation in paddy fields, and the domestication of indica in South Asia, following introgressive hybridization from the early japonica gene pool into “proto-indica” rice (about 4,000 years ago [16, 51]). At the time corresponding to the upper bound of the 95% credible interval (2,700 years ago), japonica rice and paddy field cultivation had spread to most areas of continental and insular South, East, and Southeast Asia, and indica rice was beginning to spread out of the Ganges plains (16, 52). The point estimates for the splitting of *M. oryzae* lineages correspond to the Tang Dynasty (“the Golden Age”) in China and the late classical period in India, during which food production became more rational and scientific and intensive irrigated systems of cultivation were developed, bringing about economic, demographic, and material growth (53).

Genome scans based on polymorphism and divergence revealed heterogeneity in the genomic and life history changes associated with the emergence and spread of the different lineages. Using microsatellite data and a larger collection of samples, Saleh et al. (21) identified differences in variability levels between lineages, with similar or higher

levels of genetic variability in lineages 1 and 4 than in lineages 2 and 3. Lineages 1 and 4 were also the only lineages that displayed biological features (fertile female rates and mating type ratios) consistent with sexual reproduction. Our genome-wide analyses of variability and linkage disequilibrium provided clear evidence that the continental Southeast Asian lineage 1 displays recombination and is genetically diverse, suggesting that sexual reproduction occurs and that long-term population size is relatively high, whereas pandemic lineages 2 and 3 are largely clonal and genetically depauperate, suggesting a lack of sexual reproduction and demographic bottlenecks associated with their emergence in agrosystems. However, population genomic analyses did not confirm the previously reported high variability and capacity for sexual recombination of the South Asia–United States lineage 4 (21), possibly due to differences in sample sizes between studies. The null hypothesis of clonality was not rejected by phi tests for recombination, but both total ( $\theta_w$ ) and average ( $\pi$ ) nucleotide diversity, and also the population recombination rate ( $\rho$ ), were on the same order of magnitude in lineage 4 as in lineages 2 and 3, consistent with a lack of recombination and a small effective population size.

The patterns of polymorphism and diversity at nonsynonymous and synonymous sites indicated that deleterious mutations were particularly abundant in clonal lineages 2 to 4 of *M. oryzae*, with the smaller long-term population size, consistent with a higher cost of pestification in these lineages. The introgression of genetic elements from clonal lineages harboring greater loads of deleterious mutations may counteract the efficient purging of deleterious mutations in the recombining lineage 1 from mainland Southeast Asia and lead to smaller differences in the proportion of nonsynonymous mutations between recombining and clonal lineages. However, the extensive variabilities in the origin and genomic distribution of the detected putative migrant mutations suggest that most of these mutations are false positives, with only a series of adjacent mutations of this type originating from the same donor lineage corresponding to genuine genetic exchange events. Field-scale studies in areas in which different lineages coexist should provide more detailed insights into the relative importance of interlineage recombination and make it possible to determine whether genetic exchanges are driven by positive selection or are an incidental by-product of the sympatric coexistence of interfertile lineages. We hypothesize that the accumulation of deleterious mutations in pandemic clonal complexes and gene flow into sexual lineages during disease emergence and spread are widespread phenomena, which are not due to idiosyncrasies of *M. oryzae*, and we expect these patterns to hold true in other invasive fungal plant pathogens.

An examination of additional isolates from undersampled geographic regions (including Africa and South America), based on sequencing approaches and sampling schemes tailored to detect adaptation from *de novo* mutations, will be required to enhance our understanding of the biogeography of *M. oryzae* and the genetic basis of adaptation in the different *M. oryzae* lineages. Nevertheless, the catalog of variants detected in our study provides a solid foundation for future research into the population genomics of adaptation in *M. oryzae*. Our work also provides a population-level genomic framework for defining molecular markers for the control of rice blast and investigations of the molecular basis of the differences in phenotype and fitness between divergent lineages.

## MATERIALS AND METHODS

**Genome sequencing and SNP calling.** Sequencing libraries were prepared and Illumina HiSeq 2500 sequencing was performed either at Beckman Coulter Genomics (BCG; Danvers, MA, USA) or at the Iwate Biotechnological Research Center (Table S1). Genomic DNA for sequencing at BCG was isolated from 100 mg of fresh mycelium grown in liquid medium. The mycelium was treated with enzymes degrading the cell walls (mainly beta-glucanase) and then incubated in lysis buffer (Triton 2×–1% SDS–100 mM NaCl–10 mM Tris-HCl–1 mM EDTA). Nucleic acids were extracted by treatment with chloroform:isoamyl alcohol (24:1), followed by precipitation overnight in isopropanol. They were then rinsed in 70% ethanol. The nucleic acid extract was treated with RNase A (0.2 mg/ml, final concentration) to remove RNA. The DNA was purified by another round of chloroform:isoamyl alcohol (24:1) treatment. Genomic DNA for sequencing at IBRC was isolated with a protocol adapted from the animal tissue (mouse tail) protocol



available for the Promega Wizard genomic DNA purification kit. Nucleic acids were extracted from 20 mg of fresh mycelium grown in liquid medium, which was ground into powder in liquid nitrogen with a prechilled pestle and mortar. The centrifugation time specified for the mouse tail protocol was increased to 15 min, and centrifugation was carried out at 4°C, after precipitation for 3 h at −20°C. Nucleic acids were resuspended in water, treated with RNase A (0.2 mg/ml, final concentration), purified by treatment with chloroform:isoamyl alcohol (24:1), precipitated overnight in isopropanol supplemented with 0.1 volume of sodium acetate (3 M; pH 5), and rinsed in 70% ethanol.

Sequencing reads were either paired-end reads (read length, 100 nucleotides; insert size, ~500 bp; DNAs sequenced at IBI) or single-end reads (read length, 100 nucleotides; DNAs sequenced by BCG). Reads were trimmed to remove barcodes and adapters and were then filtered to eliminate sequences containing ambiguous base calls. Reads were mapped against the 70-15 reference genome, version 8 (10), with BWA (54) (subcommand `al`, option `-n 5`; subcommand `sampe` option `-a 500`). Alignments were sorted with `samtools` (55), and reads with a mapping quality below 30 were removed. Duplicates were removed with Picard (<http://broadinstitute.github.io/picard/>). We used Realigner-Targetcreator, Targetcreator, and Indelrealigner within the genome analyses toolkit (GATK) (56) to define intervals to target for local realignment and for the local realignment of reads around indels, respectively, and Unified Genotyper to call SNPs. We used GATK's SelectVariants to apply hard filters and to select high-confidence SNPs based on annotation values. Numbers of reference and alternative alleles were calculated with JEXL expressions based on the `vc.getGenotype().getAD()` command. Variants were selected based on the following parameters: counts of all reads with a MAPQ of 0 below 3.0 (MQ0 in GATK), number of reference alleles + number of alternative alleles  $\geq 15.0$ , and number of reference alleles/number of alternative alleles  $\leq 0.1$ . With these parameters, SNP calls are limited to positions with relatively high sequencing depths and limited discordance across high-quality sequencing reads. We used a second SNP caller, Freebayes v0.9.10-3-g47a713e (57), to assess the impact of the SNP calling method on the sets of SNPs detected, given the presence in our data set of isolates sequenced at relatively low depth ( $<10\times$ ). We set the `-min-alternate-count` option to one in Freebayes. When the sample-by-sample Freebayes SNP calls were compared with the GATK SNP calls, after filtration, Freebayes identified  $1.63\times$  (standard deviation [SD], 0.28) more SNPs per sample on average than via analyses with GATK, and 92.3% (SD, 2.3) of the SNPs identified with GATK were also identified with Freebayes. The size of the intersection between the sets of SNPs identified by the two methods was negatively correlated with sequencing depth (i.e., the concordance between SNP callers was higher for isolates sequenced less deeply), indicating a minimal impact of isolates sequenced at lower depth on confidence in SNP calls. When the multisample Freebayes SNP calls were compared with the GATK SNP calls, after filtration, 83% of the SNPs identified with GATK were confirmed with Freebayes, and the GATK SNPs that were not confirmed with Freebayes were identified in sets of isolates with a genome-wide sequencing depth of  $47.8\times$  on average (SD, 8.2), consistent with a minimal impact of isolates sequenced at lower depth on confidence in SNP calls. High-confidence SNPs were annotated with SnpEff v4.3 (58).

**Mating type and female fertility assays.** Mating type and female fertility for our lineages had previously been determined (23) or we determined them as previously described (59).

**Genealogical relationships and population subdivision.** Total evidence genealogy was inferred with RAxML from pseudoassembled genomic sequences (i.e., tables of SNPs converted into a fasta file, using the reference sequence as a template), assuming a general time-reversible model of nucleotide substitution with the  $\Gamma$  model of rate heterogeneity. Bootstrap confidence levels were determined with 100 replicates. DAPC was performed with the Adegenet package in R (60). Sites with missing data were excluded. We retained the first 20 principal components and the first six discriminant functions.

**Diversity and divergence.** Polymorphism and divergence statistics were calculated with Egglib 3.0.0b10 (61), excluding sites with  $>30\%$  missing data. The neutrality index was calculated as  $(P_n/P_s)/(D_n/D_s)$ , where  $P_n$  and  $P_s$  are the numbers of nonsynonymous and synonymous polymorphisms, and  $D_n$  and  $D_s$  are the numbers of nonsynonymous and synonymous substitutions, respectively.  $D_n$  and  $D_s$  were calculated with Gestimator (62) using the *Setaria*-infecting lineage as an outgroup.  $P_n$  and  $P_s$  were calculated with Egglib.

**Linkage disequilibrium and recombination.** The coefficient of linkage disequilibrium ( $r^2$ ) (63) was calculated with Vcftools (64), excluding missing data and sites with minor allele frequencies below 10%. For all lineages, we calculated  $r^2$  values for all pairs of SNPs less than 100 kb apart and averaged LD values in distance classes of 1 kb for lineages 1 and 4 and 10 kb for lineages 2 and 3, to minimize noise due to low genetic diversity. Only sites without missing data and with a minor allele frequency above 10% were included, to minimize the dependence of  $r^2$  on minor allele frequency (65). Recombination rates were estimated for each chromosome with Pairwise in LDhat version 2.2 (66). Singletons and sites with missing data were excluded.

**Pathogenicity tests.** We used pathotyping data for 31 isolates previously described by Gallet et al. (38). We supplemented this data set with pathotyping data for 27 isolates produced by the same authors, using the same protocol but not included in their publication due to uncertainty in the nature of the rice subspecies of origin. We used a combination of multilocus microsatellite and SNP data to assign the 58 pathotyped isolates to the six lineages, because SNP data were available for only 30 pathotyped isolates (20 of the 31 isolates from Gallet et al. and 10 of the 27 additional isolates). Multilocus microsatellite genotypes at 12 loci were obtained from the Saleh et al. (21) data set or were produced as described by Saleh et al. (21). We improved the accuracy of assignment tests by adding to the full data set the 19 isolates that had been sequenced but for which no pathotyping data were available, which included 77 multilocus genotypes in total (58 pathotyped isolates and 19 additional nonpathotyped isolates). For 49



of the 77 isolates for which genomic data were available, we retained 1% of the SNP loci with no missing data (i.e., 164 SNPs). Missing data were introduced at SNP and microsatellite loci for the 28 nonsequenced isolates and the four sequenced isolates without microsatellite data.

The Structure 2.3.1 program was used for determining assignments (67–69). The model implemented allowed admixture and correlation in allele frequencies. Burn-in length was set at 10,000 iterations, and the burn-in period was followed by 40,000 iterations. Four independent runs were performed to check for convergence. At  $K = 6$ , the four main clusters identified with the full genomic data set were recovered, although 15 of the 77 genotypes could not be assigned due to admixture or a lack of power. Finally, 46 of the 58 isolates inoculated could be assigned to lineages 1 to 4; the other 12 isolates could not be assigned to a specific lineage among lineages 1, 5, and 6 and were not analyzed further (Fig. S4). Infection success was analyzed with a generalized linear model with a binomial error structure and logit link function. Treatment contrasts were used to assess the specific degrees of freedom of main effects and interactions.

**Genome scan for genetic exchanges.** Probabilistic chromosome painting was carried out with Chromopainter version 0.0.4 (70). This method “paints” individuals in “recipient” populations as a combination of segments from “donor” populations, using linkage information for probability computation and assuming that linked alleles are more likely to be exchanged together during recombination events. All lineages were used as donors, but only lineages 1 to 4 were used as recipients (sample sizes were too small for lineages 5 and 6). We initially ran the model using increments of 50 expectation-maximization iterations, starting at 10 iterations, and we examined the convergence of parameter estimates to determine how many iterations to use. Hence, the recombination scaling constant  $N_e$  and emission probabilities ( $\mu$ ) were estimated in lineages 1 to 4 by running the expectation-maximization algorithm with 200 iterations for each lineage and chromosome. Estimates of  $N_e$  and  $\mu$  were then calculated as averages weighted by chromosome length ( $N_e = 8,160$  for all lineages; lineage 1,  $\mu = 0.0000506$ ; lineage 2,  $\mu = 0.0000171$ ; lineage 3,  $\mu = 0.000021$ ; lineage 4,  $\mu = 0.000011$ ). These parameter values and the per chromosome recombination rates estimated determined with LDhat were then used to paint the chromosome of each lineage, considering the remaining lineages as donors and using 200 expectation-maximization iterations. We used a probability threshold of 0.9 to assign mutations in a recipient lineage to a donor lineage.

**Tip-calibrated phylogenetic analysis.** Tip-calibrated phylogenetic inferences were performed with only the 48 isolates for which sampling date were recorded, i.e., all isolates except the reference strain 70-15 and strain PH0018, with the exclusion of missing data. We investigated whether the signal obtained with our data set was sufficiently high for thorough tip-dating inferences by building a phylogenetic tree with PhyML (71), without constraining tip heights on the basis of isolate sampling time, and then fitting root-to-tip distances (a proxy for the number of substitutions accumulated since the most recent common ancestor [TMRCA]) to collection dates with TempEst (70). We observed a significant positive correlation (Fig. S3), demonstrating that the temporal signal was sufficiently strong for thorough tip-dating inferences at this evolutionary scale. The tip-calibrated inferences were then carried out using Markov chain-Monte Carlo sampling in beast 1.8.2 (72). The topology was fixed as the total-evidence genome genealogy inferred with RAxML. We used an annotation of the SNPs with SNPEff (57) to partition Bayesian inference (i.e., several substitution models and rates of evolution were fitted to the different sets of SNPs during a single analysis). The optimal partitioning scheme and the best-fit nucleotide substitution model for each partitioning of the genome were estimated with PartitionFinder software (73). The best partitioning was obtained for  $K = 3$  schemes (synonymous: HKY, non-synonymous: GTR and non-exonic SNPs: GTR) and was used for subsequent analyses. Node age was then estimated with this optimal partitioning scheme. Rate variation between sites was modeled with a discrete gamma distribution, with four rate categories. We assumed an uncorrelated lognormal relaxed clock, to account for rate variation between lineages. We minimized prior assumptions about demographic history, by adopting an extended Bayesian skyline plot approach, to integrate data over different coalescent histories. The tree was calibrated using tip-dates only. We applied flat priors (i.e., uniform distributions) for substitution rate ( $1 \times 10^{-12} - 1 \times 10^{-2}$  substitutions/site/year) and for the age of any internal node in the tree (including the root). We ran five independent chains, in which samples were drawn every 5,000 MCMC steps, from a total of 50,000,000 steps, after a discarded burn-in of 5,000,000 steps. We checked for convergence to the stationary distribution and for sufficient sampling and mixing by inspecting posterior samples (effective sample size,  $>200$ ). Parameter estimation was based on samples combined from the different chains. The best-supported tree was estimated from the combined samples and using the maximum clade credibility method implemented in TreeAnnotator.

**Functional enrichment.** Gene enrichment analysis was conducted with the R package TopGO for GO terms and Fisher’s exact test for enrichment in HET domain genes, NLRs, small secreted protein genes, and MAX-effector genes. MAX-effector genes were those reported by de Guillen et al. (71), NLRs were those identified by Dyrka et al. (46), and small secreted proteins and HET domain proteins were identified with Ensembl’s Biomart.

## SUPPLEMENTAL MATERIAL

Supplemental material for this article may be found at <https://doi.org/10.1128/mBio.01806-17>.

**FIG S1**, DOCX file, 0.8 MB.

**FIG S2**, DOCX file, 0.2 MB.

**FIG S3**, PDF file, 0.1 MB.

**FIG S4**, DOCX file, 0.7 MB.

**TABLE S1**, XLSX file, 0.02 MB.

**TABLE S2**, DOCX file, 0.5 MB.

**TABLE S3**, XLSX file, 0.04 MB.

**TABLE S4**, XLSX file, 0.04 MB.

**TABLE S5**, XLSX file, 0.04 MB.

**TABLE S6**, XLSX file, 0.04 MB.

## ACKNOWLEDGMENTS

We thank the scholars who contributed samples, François Bonnot and Romain Gallet for assistance with statistics, and the South Green and Migale computing facilities.

## REFERENCES

- Taylor JW, Jacobson D, Fisher M. 1999. The evolution of asexual fungi: reproduction, speciation and classification. *Annu Rev Phytopathol* 37: 197–246. <https://doi.org/10.1146/annurev.phyto.37.1.197>.
- Gladieux P, Feurtey A, Hood ME, Snirc A, Clavel J, Dutech C, Roy M, Giraud T. 2015. The population biology of fungal invasions. *Mol Ecol* 24:1969–1986. <https://doi.org/10.1111/mec.13028>.
- Martin MD, Vieira FG, Ho SYW, Wales N, Schubert M, Seguin-Orlando A, Ristaino JB, Gilbert MT. 2016. Genomic characterization of a South American *Phytophthora* hybrid mandates reassessment of the geographic origins of *Phytophthora infestans*. *Mol Biol Evol* 33:478–491. <https://doi.org/10.1093/molbev/msv241>.
- Simwami SP, Khayhan K, Henk DA, Aanensen DM, Boekhout T, Hagen F, Brouwer AE, Harrison TS, Donnelly CA, Fisher MC. 2011. Low diversity *Cryptococcus neoformans* variety *grubii* multilocus sequence types from Thailand are consistent with an ancestral African origin. *PLoS Pathog* 7:e1001343. <https://doi.org/10.1371/journal.ppat.1001343>.
- Ali S, Gladieux P, Rahman H, Saqib MS, Fiaz M, Ahmad H, Leconte M, Gautier A, Justesen AF, Hovmøller MS, Enjalbert J, de Vallavieille-Pope C. 2014. Inferring the contribution of sexual reproduction, migration and off-season survival to the temporal maintenance of microbial populations: a case study on the wheat fungal pathogen *Puccinia striiformis* f. sp. *tritici*. *Mol Ecol* 23:603–617. <https://doi.org/10.1111/mec.12629>.
- Savary S, Willocquet L, Elazegui FA, Castilla NP, Teng PS. 2000. Rice pest constraints in tropical Asia: quantification of yield losses due to rice pests in a range of production situations. *Plant Dis* 84:357–369. <https://doi.org/10.1094/PDIS.2000.84.3.357>.
- Talbot NJ. 2003. On the trail of a cereal killer: exploring the biology of *Magnaporthe grisea*. *Annu Rev Microbiol* 57:177–202. <https://doi.org/10.1146/annurev.micro.57.030502.090957>.
- Gurr S, Samalova M, Fisher M. 2011. The rise and rise of emerging infectious fungi challenges food security and ecosystem health. *Fungal Biol Rev* 25:181–188. <https://doi.org/10.1016/j.fbr.2011.10.004>.
- Valent B. 1990. Rice blast as a model system for plant pathology. *Phytopathology* 80:33–36. <https://doi.org/10.1094/phyto-80-33>.
- Dean RA, Talbot NJ, Ebbole DJ, Farman ML, Mitchell TK, Orbach MJ, Thon M, Kulkarni R, Xu JR, Pan H, Read ND, Lee YH, Carbone I, Brown D, Oh YY, Donofrio N, Jeong JS, Soanes DM, Djonovic S, Kolomiets E, Rehmeier C, Li W, Harding M, Kim S, Lebrun MH, Bohnert H, Coughlan S, Butler J, Calvo S, Ma LJ, Nicol R, Purcell S, Nusbaum C, Galagan JE, Birren BW. 2005. The genome sequence of the rice blast fungus *Magnaporthe grisea*. *Nature* 434:980–986. <https://doi.org/10.1038/nature03449>.
- Ebbole DJ. 2007. *Magnaporthe* as a model for understanding host-pathogen interactions. *Annu Rev Phytopathol* 45:437–456. <https://doi.org/10.1146/annurev.phyto.45.062806.094346>.
- Chiapello H, Mallet L, Guérin C, Aguilera G, Amselem J, Kroj T, Ortega-Abboud E, Lebrun MH, Henrissat B, Gendral A, Rodolphe F, Tharreau D, Fournier E. 2015. Deciphering genome content and evolutionary relationships of isolates from the fungus *Magnaporthe oryzae* attacking different host plants. *Genome Biol Evol* 7:2896–2912. <https://doi.org/10.1093/gbe/evv187>.
- Islam MT, Croll D, Gladieux P, Soanes DM, Persoons A, Bhattacharjee P, Hossain MS, Gupta DR, Rahman MM, Mahboob MG, Cook N, Salam MU, Surovy MZ, Sancho VB, Maciel JL, Nhani Júnior A, Castroagudín VL, Reges JT, Ceresini PC, Ravel S, Kellner R, Fournier E, Tharreau D, Lebrun MH, McDonald BA, Stitt T, Swan D, Talbot NJ, Saunders DG, Win J, Kamoun S. 2016. Emergence of wheat blast in Bangladesh was caused by a South American lineage of *Magnaporthe oryzae*. *BMC Biol* 14:84. <https://doi.org/10.1186/s12915-016-0309-7>.
- Yoshida K, Saunders DGO, Mitsuoka C, Natsume S, Kosugi S, Saitoh H, Inoue Y, Chuma I, Tosa Y, Cano LM, Kamoun S, Terauchi R. 2016. Host specialization of the blast fungus *Magnaporthe oryzae* is associated with dynamic gain and loss of genes linked to transposable elements. *BMC Genomics* 17:370. <https://doi.org/10.1186/s12864-016-2690-6>.
- Couch BC, Fudal I, Lebrun MH, Tharreau D, Valent B, van Kim P, Nottéghem JL, Kohn LM. 2005. Origins of host-specific populations of the blast pathogen *Magnaporthe oryzae* in crop domestication with subsequent expansion of pandemic clones on rice and weeds of rice. *Genetics* 170:613–630. <https://doi.org/10.1534/genetics.105.041780>.
- Fuller DQ, Sato Y-I, Castillo C, Qin L, Weisskopf AR, Kingwell-Banham EJ, Song J, Ahn S, van Etten J. 2010. Consilience of genetics and archaeobotany in the entangled history of rice. *Archaeol Anthropol Sci* 2:115–131. <https://doi.org/10.1007/s12520-010-0035-y>.
- Diao X, Jia G. 2017. Origin and domestication of foxtail millet, p 61–72. In Doust A, Diao X (ed), *Genetics and genomics of Setaria*. Springer, Cham, Switzerland.
- Huang X, Kurata N, Wei X, Wang ZX, Wang A, Zhao Q, Zhao Y, Liu K, Lu H, Li W, Guo Y, Lu Y, Zhou C, Fan D, Weng Q, Zhu C, Huang T, Zhang L, Wang Y, Feng L, Furuumi H, Kubo T, Miyabayashi T, Yuan X, Xu Q, Dong G, Zhan Q, Li C, Fujiyama A, Toyoda A, Lu T, Feng Q, Qian Q, Li J, Han B. 2012. A map of rice genome variation reveals the origin of cultivated rice. *Nature* 490:497–501. <https://doi.org/10.1038/nature11532>.
- Huang X, Han B. 2015. Rice domestication occurred through single origin and multiple introgressions. *Nat Plants* 2:15207. <https://doi.org/10.1038/nplants.2015.207>.
- Choi JY, Platts AE, Fuller DQ, Hsing YI, Wing RA, Purugganan MD. 2017. The rice paradox: multiple origins but single domestication in Asian rice. *Mol Biol Evol* 34:969–979. <https://doi.org/10.1093/molbev/msx049>.
- Saleh D, Milazzo J, Adreit H, Fournier E, Tharreau D. 2014. South-East Asia is the center of origin, diversity and dispersion of the rice blast fungus, *Magnaporthe oryzae*. *New Phytol* 201:1440–1456. <https://doi.org/10.1111/nph.12627>.
- Zeigler RS. 1998. Recombination in *Magnaporthe grisea*. *Annu Rev Phytopathol* 36:249–275. <https://doi.org/10.1146/annurev.phyto.36.1.249>.
- Saleh D, Xu P, Shen Y, Li C, Adreit H, Milazzo J, Ravigné V, Bazin E, Nottéghem JL, Fournier E, Tharreau D. 2012. Sex at the origin: an Asian population of the rice blast fungus *Magnaporthe oryzae* reproduces sexually. *Mol Ecol* 21:1330–1344. <https://doi.org/10.1111/j.1365-294X.2012.05469.x>.
- Lu J, Tang T, Tang H, Huang J, Shi S, Wu CI. 2006. The accumulation of deleterious mutations in rice genomes: a hypothesis on the cost of domestication. *Trends Genet* 22:126–131. <https://doi.org/10.1016/j.tig.2006.01.004>.
- Glémin S, Bataillon T. 2009. A comparative view of the evolution of grasses under domestication. *New Phytol* 183:273–290. <https://doi.org/10.1111/j.1469-8137.2009.02884.x>.
- Stukenbrock EH, Bataillon T, Dutheil JY, Hansen TT, Li R, Zala M, McDonald BA, Wang J, Schierup MH. 2011. The making of a new pathogen:

- insights from comparative population genomics of the domesticated wheat pathogen *Mycosphaerella graminicola* and its wild sister species. *Genome Res* 21:2157–2166. <https://doi.org/10.1101/gr.118851.110>.
27. Stukenbrock EH, Bataillon T. 2012. A population genomics perspective on the emergence and adaptation of new plant pathogens in agroecosystems. *PLoS Pathog* 8:e1002893. <https://doi.org/10.1371/journal.ppat.1002893>.
  28. Felsenstein J. 1974. The evolutionary advantage of recombination. *Genetics* 78:737–756.
  29. Karasov T, Messer PW, Petrov DA. 2010. Evidence that adaptation in *Drosophila* is not limited by mutation at single sites. *PLoS Genet* 6:e1000924. <https://doi.org/10.1371/journal.pgen.1000924>.
  30. Ellison CE, Hall C, Kowbel D, Welch J, Brem RB, Glass NL, Taylor JW. 2011. Population genomics and local adaptation in wild isolates of a model microbial eukaryote. *Proc Natl Acad Sci U S A* 108:2831–2836. <https://doi.org/10.1073/pnas.1014971108>.
  31. Roper M, Ellison C, Taylor JW, Glass NL. 2011. Nuclear and genome dynamics in multinucleate ascomycete fungi. *Curr Biol* 21:R786–R793. <https://doi.org/10.1016/j.cub.2011.06.042>.
  32. Cheeseman K, Ropars J, Renault P, Dupont J, Gouzy J, Branca A, Abraham AL, Ceppi M, Conseiller E, Debuchy R, Malagnac F, Goarin A, Silar P, Lacoste S, Sallet E, Bensimon A, Giraud T, Brygoo Y. 2014. Multiple recent horizontal transfers of a large genomic region in cheese making fungi. *Nat Commun* 5:2876. <https://doi.org/10.1038/ncomms3876>.
  33. Gladieux P, Ropars J, Badouin H, Branca A, Aguilera G, De Vienne DM, Rodríguez de la Vega RC, Branco S, Giraud T. 2014. Fungal evolutionary genomics provides insight into the mechanisms of adaptive divergence in eukaryotes. *Mol Ecol* 23:753–773. <https://doi.org/10.1111/mec.12631>.
  34. Noguchi MT, Yasuda N, Fujita Y. 2006. Evidence of genetic exchange by parasexual recombination and genetic analysis of pathogenicity and mating type of parasexual recombinants in rice blast fungus, *Magnaporthe oryzae*. *Phytopathology* 96:746–750. <https://doi.org/10.1094/PHYTO-96-0746>.
  35. Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* 30:1312–1313. <https://doi.org/10.1093/bioinformatics/btu033>.
  36. Leaché AD, Banbury BL, Felsenstein J, de Oca AN, Stamatakis A. 2015. Short tree, long tree, right tree, wrong tree: new acquisition bias corrections for inferring SNP phylogenies. *Syst Biol* 64:1032–1047. <https://doi.org/10.1093/sysbio/syv053>.
  37. Akashi H, Osada N, Ohta T. 2012. Weak selection and protein evolution. *Genetics* 192:15–31. <https://doi.org/10.1534/genetics.112.140178>.
  38. Gallet R, Fontaine C, Bonnot F, Milazzo J, Tertois C, Adreit H, Ravigné V, Fournier E, Tharreau D. 2016. Evolution of compatibility range in the rice-Magnaporthe oryzae system: an uneven distribution of R genes between rice subspecies. *Phytopathology* 106:348–354. <https://doi.org/10.1094/PHYTO-07-15-0169-R>.
  39. Giraud T, Gladieux P, Gavrillets S. 2010. Linking emergence of fungal plant diseases with ecological speciation. *Trends Ecol Evol* 25:387–395. <https://doi.org/10.1016/j.tree.2010.03.006>.
  40. Schulze-Lefert P, Panstruga R. 2011. A molecular evolutionary concept connecting nonhost resistance, pathogen host range, and pathogen speciation. *Trends Plant Sci* 16:117–125. <https://doi.org/10.1016/j.tplants.2011.01.001>.
  41. Liao J, Huang H, Meusnier I, Adreit H, Ducasse A, Bonnot F, Pan L, He X, Kroj T, Fournier E, Tharreau D, Gladieux P, Morel JB. 2016. Pathogen effectors and plant immunity determine specialization of the blast fungus to rice subspecies. *eLife* 5:e19377. <https://doi.org/10.7554/eLife.19377>.
  42. Bryant D, Moulton V. 2004. Neighbor-Net: an agglomerative method for the construction of phylogenetic networks. *Mol Biol Evol* 21:255–265. <https://doi.org/10.1093/molbev/msh018>.
  43. Bruen TC, Philippe H, Bryant D. 2006. A simple and robust statistical test for detecting the presence of recombination. *Genetics* 172:2665–2681. <https://doi.org/10.1534/genetics.105.048975>.
  44. Maddison WP, Maddison DR. 2016. Mesquite: a modular system for evolutionary analysis, version 3.11. <http://mesquiteproject.org>.
  45. Ardlie KG, Seielstad M. 2002. Patterns of linkage disequilibrium in the human genome. *Nat Rev Genet* 3:299–309. <https://doi.org/10.1038/nrg777>.
  46. Dyrka W, Lamacchia M, Durrens P, Kobe B, Daskalov A, Paoletti M, Sherman DJ, Saupe SJ. 2014. Diversity and variability of NOD-like receptors in fungi. *Genome Biol Evol* 6:3137–3158. <https://doi.org/10.1093/gbe/evu251>.
  47. Saupe SJ. 2000. Molecular genetics of heterokaryon incompatibility in filamentous ascomycetes. *Microbiol Mol Biol Rev* 64:489–502. <https://doi.org/10.1128/MMBR.64.3.489-502.2000>.
  48. Rieux A, Balloux F. 2016. Inferences from tip-calibrated phylogenies: a review and a practical guide. *Mol Ecol* 25:1911–1924. <https://doi.org/10.1111/mec.13586>.
  49. Zhao Z. 2011. New archaeobotanic data for the study of the origins of agriculture in China. *Curr Anthropol* 52:S295–S306. <https://doi.org/10.1086/659308>.
  50. Castillo CC, Bellina B, Fuller DQ. 2016. Rice, beans and trade crops on the early maritime Silk Route in Southeast Asia. *Antiquity* 90:1255–1269. <https://doi.org/10.15184/aqy.2016.175>.
  51. Fujiwara H. 1996. Search for the origin of rice cultivation: the ancient rice cultivation in paddy fields at the Cao Xie Shan Site in China. Society for Scientific Studies on Cultural Property, Miyazaki, Japan.
  52. Fuller DQ, Qin L. 2009. Water management and labour in the origins and dispersal of Asian rice. *World Archaeol* 41:88–111. <https://doi.org/10.1080/00438240802668321>.
  53. Siddiqui IH. 2008. Water works and irrigation system in India during pre-Mughal times, p 429–454. Brill Online, Leiden, The Netherlands.
  54. Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25:1754–1760. <https://doi.org/10.1093/bioinformatics/btp324>.
  55. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, 1000 Genome Project Data Processing Subgroup. 2009. The sequence alignment/map format and SAMtools. *Bioinformatics* 25:2078–2079. <https://doi.org/10.1093/bioinformatics/btp352>.
  56. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A, Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA. 2010. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res* 20:1297–1303. <https://doi.org/10.1101/gr.107524.110>.
  57. Garrison E, Marth G. 2012. Haplotype-based variant detection from short-read sequencing. *arXiv arXiv:12073907*. [q-bio.GN]. <https://arxiv.org/abs/1207.3907>.
  58. Cingolani P, Platts A, Wang le L, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM. 2012. A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w<sup>1118</sup>;iso-2; iso-3. *Fly* 6:80–92. <https://doi.org/10.4161/fly.19695>.
  59. Notteghem JL, Silue D. 1992. Distribution of the mating type alleles in *Magnaporthe grisea* populations pathogenic on rice. *Phytopathology* 82:421–424. <https://doi.org/10.1094/Phyto-82-421>.
  60. Jombart T, Ahmed I. 2011. Adegenet 1.3-1: new tools for the analysis of genome-wide SNP data. *Bioinformatics* 27:3070–3071. <https://doi.org/10.1093/bioinformatics/btr521>.
  61. De Mita S, Siol M. 2012. EggLib: processing, analysis and simulation tools for population genetics and genomics. *BMC Genet* 13:27. <https://doi.org/10.1186/1471-2156-13-27>.
  62. Thornton K. 2003. Libsequence: a C++ class library for evolutionary genetic analysis. *Bioinformatics* 19:2325–2327. <https://doi.org/10.1093/bioinformatics/btg316>.
  63. Hill WG, Robertson A. 1968. Linkage disequilibrium in finite populations. *Theor Appl Genet* 38:226–231. <https://doi.org/10.1007/BF01245622>.
  64. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE, Lunter G, Marth GT, Sherry ST, McVean G, Durbin R, 1000 Genomes Project Analysis Group. 2011. The variant call format and VCFtools. *Bioinformatics* 27:2156–2158. <https://doi.org/10.1093/bioinformatics/btr330>.
  65. Amaral AJ, Megens HJ, Crooijmans RPMA, Heuven HCM, Groenen MAM. 2008. Linkage disequilibrium decay and haplotype block structure in the pig. *Genetics* 179:569–579. <https://doi.org/10.1534/genetics.107.084277>.
  66. Auton A, McVean G. 2007. Recombination rate estimation in the presence of hotspots. *Genome Res* 17:1219–1227. <https://doi.org/10.1101/gr.6386707>.
  67. Pritchard JK, Stephens M, Donnelly P. 2000. Inference of population structure using multilocus genotype data. *Genetics* 155:945–959.
  68. Falush D, Stephens M, Pritchard JK. 2003. Inference of population structure using multilocus genotype data: linked loci and correlated allele frequencies. *Genetics* 164:1567–1587.
  69. Hubisz MJ, Falush D, Stephens M, Pritchard JK. 2009. Inferring weak population structure with the assistance of sample group information. *Mol Ecol Resour* 9:1322–1332. <https://doi.org/10.1111/j.1755-0998.2009.02591.x>.
  70. Lawson DJ, Hellenthal G, Myers S, Falush D. 2012. Inference of popula-

- tion structure using dense haplotype data. *PLoS Genet* 8:e1002453. <https://doi.org/10.1371/journal.pgen.1002453>.
71. de Guillen K, Ortiz-Vallejo D, Gracy J, Fournier E, Kroj T, Padilla A. 2015. Structure analysis uncovers a highly diverse but structurally conserved effector family in phytopathogenic fungi. *PLoS Pathog* 11:e1005228. <https://doi.org/10.1371/journal.ppat.1005228>.
72. Drummond AJ, Suchard MA, Xie D, Rambaut A. 2012. Bayesian phylogenetics with BEAUti and the BEAST 1.7. *Mol Biol Evol* 29:1969–1973.
73. Lanfear R, Frandsen PB, Wright AM, Senfeld T, Calcott B. 2016. Partition-Finder 2: new methods for selecting partitioned models of evolution for molecular and morphological phylogenetic analyses. *Mol Biol Evol* 34: 772–773. <https://doi.org/10.1093/molbev/msw260>.